



# Digital Twin-Enabled Multigeneration Control Co-Design With Deep Reinforcement Learning

**Ying-Kuan Tsai**

Department of Mechanical Engineering,  
 Northwestern University,  
 Evanston, IL 60208  
 e-mail: [yingkuan.tsai@northwestern.edu](mailto:yingkuan.tsai@northwestern.edu)

**Vispi Karkaria**

Department of Mechanical Engineering,  
 Northwestern University,  
 Evanston, IL 60208  
 e-mail: [vispikarkaria2026@u.northwestern.edu](mailto:vispikarkaria2026@u.northwestern.edu)

**Yi-Ping Chen**

Department of Mechanical Engineering,  
 Northwestern University,  
 Evanston, IL 60208  
 e-mails: [chenyp@u.northwestern.edu](mailto:chenyp@u.northwestern.edu);  
[yi-pingchen2026@u.northwestern.edu](mailto:yi-pingchen2026@u.northwestern.edu)

**Wei Chen<sup>1</sup>**

Department of Mechanical Engineering,  
 Northwestern University,  
 Evanston, IL 60208  
 e-mail: [weichen@northwestern.edu](mailto:weichen@northwestern.edu)

*Control co-design (CCD) integrates physical and control system design to improve the performance of dynamic and autonomous systems. Despite advances in uncertainty-aware CCD methods, real-world uncertainties remain highly unpredictable. Multigeneration design addresses this challenge by considering the full lifecycle of a product: data collected from each generation informs the design of subsequent generations, enabling progressive improvements in robustness and efficiency. Digital twin (DT) technology further strengthens this paradigm by creating virtual representations that evolve over the lifecycle through real-time sensing, model updating, and adaptive re-optimization. This article presents a DT-enabled CCD framework that integrates deep reinforcement learning (DRL) to jointly optimize physical design and controller. DRL accelerates real-time decision-making by allowing controllers to continuously learn from data and adapt to uncertain environments. Extending this approach, the framework employs a multigeneration paradigm, where each cycle of deployment, operation, and redesign uses collected data to refine DT models, improve uncertainty quantification through quantile regression, and inform next-generation designs of both physical components and controllers. The framework is demonstrated on an active suspension system, where DT-enabled learning from road conditions and driving behaviors yields smoother and more stable control trajectories. Results show that the method significantly enhances dynamic performance, robustness, and efficiency. Contributions of this work include: (1) extending CCD into a lifecycle-oriented multigeneration framework, (2) leveraging DTs for continuous model updating and informed design, and (3) employing DRL to accelerate adaptive real-time decision-making. [DOI: 10.1115/1.4072056]*

*Keywords: digital twin, control co-design, deep reinforcement learning, uncertainty quantification, multigeneration design, quantile regression, artificial intelligence, data-driven design, machine learning, systems design*

## 1 Introduction

Autonomous systems are increasingly being adopted across various domains, from self-driving vehicles to robotics, where adaptability to uncertain environments is critical. Traditional design approaches often treat the physical system design and control policy separately, leading to suboptimal solutions. Control co-design (CCD) addresses this limitation by simultaneously optimizing both the physical system and its controller to improve performance [1–3]. For example, active suspension systems in vehicles offer significant advantages over passive suspension systems, which rely solely on mechanical properties to absorb road disturbances [4–7]. Active suspension employs real-

time control adjustments to enhance ride comfort, stability, and safety. By leveraging CCD, the design and control of active suspension can be jointly optimized, leading to superior performance in diverse and uncertain driving conditions.

However, achieving robust CCD in real-world applications presents several challenges. Vehicle suspension systems, for example, operate under highly unpredictable conditions, such as varying road surfaces, changing weather patterns, and diverse driver behaviors. Some existing methods focus on uncertainty-aware CCD to address these uncertainties, where they are explicitly considered during the design stage to generate robust solutions, including stochastic programming [8], robust model predictive control (MPC) [9,10], stochastic MPC [6,11], reliability-based design optimization [12], and robust design optimization [13,14]. Recent efforts have also explored real-time adaptability and multifidelity surrogate modeling within CCD frameworks [15,16].

While these methods improve robustness and computational efficiency with the consideration of uncertainty, they typically assume well-defined disturbances and may struggle to adapt to

<sup>1</sup>Corresponding author.

Contributed by the Design Automation Committee of ASME for publication in the JOURNAL OF MECHANICAL DESIGN. Manuscript received September 22, 2025; final manuscript received May 18, 2026; published online June 17, 2026. Assoc. Editor: Douglas Allaire.

unforeseen variations in real-world conditions. Additionally, these model-based approaches present a critical challenge in real-world applications, limiting their practical applicability. These limitations highlight the need for approaches that can go beyond predefined disturbance models and generalize across diverse operating conditions, such as estimating uncertain parameters [17] or capturing unmodeled dynamical behaviors [18]. Data-driven methods offer a promising alternative for solving CCD problems by leveraging real-world data to continuously update system models, enable data-driven characterization of model discrepancy and predictive uncertainty from deployment data, and support improved adaptation across evolving operating conditions in dynamic environments.

Numerous studies have integrated reinforcement learning (RL) techniques into CCD, enabling autonomous systems to cooptimize both their physical design and control policies based on performance feedback [19]. Unlike traditional model-based methods, RL-based CCD frameworks jointly optimize both physical structures and control policies by learning directly from interactions with the environment, making them well-suited for handling high-dimensional design spaces and adapting to changing conditions. In contrast to model-based approaches that rely on explicitly defined disturbance models, learning-based CCD frameworks treat the environment as partially unknown and allow the agent to learn control and design policies through interaction with the system. Although disturbances in the simulation environment are generated from predefined distributions for evaluation purposes, the learning agent does not have prior knowledge of these disturbances and must adapt its policy based on observed system responses. Two major formulations have been explored. The first is bi-level optimization, where the high-level process searches for optimal hardware parameters while the low-level RL agent learns an optimal control policy through trial and error [20]. To improve sample efficiency and exploration, techniques such as Bayesian optimization and evolutionary algorithms have been integrated into this framework [21,22].

A second approach is the “hardware as policy” paradigm, where the mechanical design itself is embedded into the control policy and jointly optimized with the parameters of the RL controller [23,24]. By representing physical parameters as differentiable components in a computational graph, gradient-based optimization algorithms can efficiently cooptimize neural network weights and mechanical design variables simultaneously [23–26]. This enables more efficient and coordinated adaptation compared to bi-level methods, which treat design and control separately and often suffer from slow convergence. Such simultaneous optimization mirrors natural coevolution, where animals adapt their body structures and behaviors together in response to growth and changing environments. While these methods have demonstrated success in enhancing adaptability, they often assume fixed training environments and struggle to incorporate real-world uncertainties. This potentially limits their robustness in practical applications.

The recent advancement of digital twin (DT) technology offers a transformative solution to CCD and other optimization problems [27]. By creating virtual representations of physical systems that evolve alongside their real-world counterparts, DTs facilitate continuous data assimilation and adaptive re-optimization [28–34]. This capability is particularly valuable for dynamic systems, such as vehicle suspension design, where real-time operational data can be leveraged to personalize suspension settings based on individual driving styles and preferences. Unlike traditional CCD approaches that rely on predefined models and assumptions, DTs provide a data-driven framework that continuously refines system models based on real-world conditions, making them more adaptable to varying operational demands.

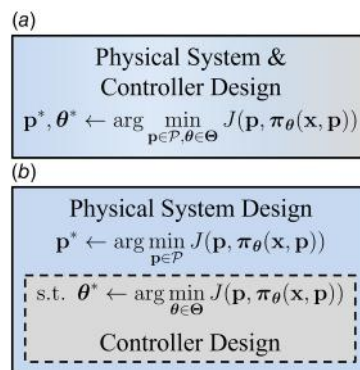
Beyond real-time optimization, multigenerational design plays a crucial role in the evolution of engineering systems [33]. While real-time updates of digital models and control policies by DTs can enhance adaptability, reoptimization of physical components using accumulated data enables long-term design improvements, refining system configurations and material choices beyond

immediate operational needs. From one generation to the next, designers and engineers leverage data collected from the current generation to inform decision-making and guide system redesign [33]. This approach ensures that each successive iteration benefits from accumulated operational insights to improve performance, reliability, and efficiency. In the context of vehicle dynamics, DTs can simulate different road conditions, driving behaviors, and suspension configurations, allowing engineers to iteratively improve the design and operation of suspension systems [35]. When the system is deployed in real-world conditions, RL controllers are continuously adapted and optimized based on online feedback data from sensors [36].

While previous works on RL-based DTs have demonstrated their effectiveness in optimizing system performance [37] through the change in control policies, none have fully explored the simultaneous optimization of both physical system design and control policies. This article addresses this gap by proposing a CCD framework for DT-enabled systems using deep reinforcement learning (DRL). Our prior work has applied reinforcement learning-based co-design to full-vehicle active suspension systems with higher-dimensional dynamics and partial observability, demonstrating the practical applicability of RL-based CCD in complex engineering systems [38]. The DT framework enables continual model refinement through data collected during operation, allowing the system to progressively improve its performance under previously unseen conditions. In addition, multigeneration design concepts are introduced, where data from previous generations are leveraged to improve both the accuracy of the digital model and system performance for future generations of physical systems. This research makes the following key contributions:

- We extend CCD to multigenerational design following the DT framework that enables the physical system and control policy to coadapt to dynamic and varying environments and improve overall performance over multiple generations of physical systems.
- We propose a learning-based CCD method by integrating DRL with data-driven models to achieve rapid real-time decision-making and support iterative system redesign with model updates.
- We incorporate uncertainty quantification (UQ) methods using physical data from previous generations to help designers make informed and robust decisions for next-generation designs.

The remainder of this article is organized as follows: Sec. 2 presents the preliminaries and technical background; Sec. 3 presents the proposed CCD method and framework for DT-enabled systems with RL; Sec. 4 implements and demonstrates the



**Fig. 1 CCD formulations for feedback policy-based systems: (a) simultaneous and (b) nested (bi-level), where  $\pi_\theta$  represents a feedback policy, which is a function of states  $\mathbf{x}$  and system parameters  $\mathbf{p} \in \mathcal{P}$  with policy parameters  $\theta \in \Theta$**

proposed approach with an engineering study of an active suspension system; and Sec. 5 summarizes the conclusion and the future work.

## 2 Technical Background

**Notation:** The sets of real numbers and nonnegative integers are denoted by  $\mathbb{R}$  and  $\mathbb{N}_{\geq 0}$ , respectively. Given  $a, b \in \mathbb{N}_{\geq 0}$  such that  $a < b$ , we denote  $\mathbb{N}_{[a,b]} := \{a, a+1, \dots, b\}$ .  $\mathbf{x}_k$  denotes the state  $\mathbf{x}$  at time  $k$ . The notation  $\mathbf{I}_{a \times a}$  denotes an  $a$ -by- $a$  identity matrix. Given a random variable  $X$ ,  $\mathbb{E}[X]$  denotes its expected value. A Gaussian distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$  is represented as  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .

**2.1 Control Co-Design.** CCD formulations integrate the physical system and controller design problems to explicitly account for the coupling between system dynamics and control strategies. This approach helps designers achieve system-level optimum compared to traditional sequential design methods [1–3]. There are two most common CCD formulations (simultaneous and nested) [39]. Figure 1 compares the simultaneous and nested formulations with a feedback control policy  $\boldsymbol{\pi}_\theta$  where  $\theta$  denotes the vector of policy parameters. Simultaneous CCD optimizes both physical and control design variables within a single optimization problem, while nested CCD separates physical and control design into an outer-loop (physical design) and an inner-loop (control optimization). Nested CCD is useful when the control problem requires specialized techniques or when different objective functions govern the physical and control system designs [9–11].

The selection of CCD formulations depends on problem complexity, computational efficiency, and the availability of specialized control design tools. In this article, we propose a CCD approach that uses the simultaneous CCD formulation to synergistically optimize both physical system variables and control policies using learning-based techniques. The proposed method captures the intricate interactions between system design and control, allowing the system and the controller to coadapt to dynamic and uncertain environments by learning from the real-time collected data.

**2.2 Reinforcement Learning.** RL is a framework for sequential decision-making where an agent interacts with an environment to learn an optimal policy that maximizes a cumulative reward [40], shown as Fig. 2. The RL problem is typically formulated as a Markov decision process (MDP), defined by a tuple  $(\mathcal{X}, \mathcal{U}, P, R)$ , where  $\mathcal{X}$  is state space,  $\mathcal{U}$  is action (control input) space,  $P$  is transition function (dynamics) of a system, and  $R$  is reward function (representing negative costs). In an MDP, the state transitions are typically stochastic, meaning that the next state  $\mathbf{x}_{k+1}$  is not deterministic when an action  $\mathbf{u}_k$  is applied to the system at state  $\mathbf{x}_k$  at time  $k \in \mathbb{N}_{\geq 0}$ . The transition function  $P(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}_k)$  represents the probability of reaching state  $\mathbf{x}_{k+1}$  from state  $\mathbf{x}_k$  by taking action  $\mathbf{u}_k$ . However, in dynamics and controls, we prefer to represent the dynamic behavior by expressing the next state as a function of the current state and control action:

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) \quad (1)$$

which is a stochastic process according to the transition function  $P(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}_k)$ .

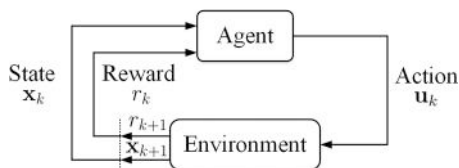


Fig. 2 Diagram of RL, modified from Ref. [41]

Once the system transits to  $\mathbf{x}_{k+1}$ , a scalar reward  $r_{k+1} = R(\mathbf{x}_k, \mathbf{u}_k, \mathbf{x}_{k+1})$  is obtained, based on the reward function  $R: \mathcal{X} \times \mathcal{U} \times \mathcal{X} \mapsto \mathbb{R}$ . A policy is a mapping from states to probabilities of selecting each possible action. That is, the agent of RL chooses action  $\mathbf{u}_k$  given state  $\mathbf{x}_k$  is  $\boldsymbol{\pi}(\mathbf{u}_k|\mathbf{x}_k)$ . Given  $(\mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1, \dots)$  with the associated rewards  $r_1, r_2, \dots$ , the return can be expressed as

$$G_k = \sum_{i=0}^{\infty} \gamma^i r_{k+i+1} \quad (2)$$

where  $\gamma \in (0, 1]$  is the discount factor. The value function of state  $\mathbf{x}$  under a policy  $\boldsymbol{\pi}$ , denoted as

$$\begin{aligned} V^\pi(\mathbf{x}) &= \mathbb{E}_\pi[G_k | \mathbf{x}_k = \mathbf{x}] \\ &= \mathbb{E}_\pi \left[ \sum_{i=0}^{\infty} \gamma^i R(\mathbf{x}_{k+i}, \mathbf{u}_{k+i}, \mathbf{x}_{k+i+1}) \mid \mathbf{x}_k = \mathbf{x} \right] \end{aligned} \quad (3)$$

for all  $\mathbf{x} \in \mathcal{X}$  represents the expected return when starting in state  $\mathbf{x}$  and following  $\boldsymbol{\pi}$  thereafter. We also call the function  $V^\pi(\mathbf{x})$  state-value function for policy  $\boldsymbol{\pi}$  or V-function.

The state-value function  $V^\pi(\mathbf{x})$  only tells us the expected return from a state  $\mathbf{x}$ , assuming that the policy  $\boldsymbol{\pi}$  is followed. However, when making decisions, we need to evaluate specific actions rather than just states. We define the value of taking action  $\mathbf{u}$  in state  $\mathbf{x}$  under a policy  $\boldsymbol{\pi}$ , denoted

$$\begin{aligned} Q^\pi(\mathbf{x}, \mathbf{u}) &= \mathbb{E}_\pi[G_k | \mathbf{x}_k = \mathbf{x}, \mathbf{u}_k = \mathbf{u}] \\ &= \mathbb{E}_\pi \left[ \sum_{i=0}^{\infty} \gamma^i R(\mathbf{x}_{k+i}, \mathbf{u}_{k+i}, \mathbf{x}_{k+i+1}) \mid \mathbf{x}_k = \mathbf{x}, \mathbf{u}_k = \mathbf{u} \right] \end{aligned} \quad (4)$$

for all  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{u} \in \mathcal{U}$ . We also call the function  $Q^\pi(\mathbf{x}, \mathbf{u})$  action-value function for policy  $\boldsymbol{\pi}$  or Q-function.  $Q(\mathbf{x}, \mathbf{u})$  provides this information by estimating the expected return for taking action  $\mathbf{u}$  in state  $\mathbf{x}$  and then following policy  $\boldsymbol{\pi}$ .

Traditional RL methods, such as tabular Q-learning and dynamic programming, struggle when applied to complex control problems with continuous state and action spaces. In such cases, it becomes nearly impossible to explicitly store or update the value function for every possible state–action pair due to the infinite number of possibilities. This limitation makes traditional RL methods impractical for high-dimensional control problems, such as robotic manipulation, autonomous driving, or real-time optimization in DT-enabled systems.

DRL addresses this challenge by leveraging deep neural networks (DNNs) to approximate value functions, policies, or both [40,42]. Instead of maintaining a lookup table for state–action values, DRL methods use neural networks as function approximators to generalize across similar states and actions. This enables RL algorithms to handle continuous, high-dimensional spaces efficiently. In particular, DRL has shown remarkable success in domains such as robotics, autonomous driving, and game playing, where large state–action spaces make traditional methods impractical. For example, policy-based methods, like trust region policy optimization [43] and proximal policy optimization (PPO) [44], are widely used DRL algorithms because the algorithms exhibit more efficient and stable policy learning by incorporating trust region optimization. This ensures that policy updates remain within a constrained step size, thus preventing drastic changes that could destabilize learning. In addition to the aforementioned methods, numerous other DRL algorithms have been developed, each tailored to specific challenges such as sample efficiency, exploration, and stability. For a comprehensive review of DRL techniques, readers are referred to Refs. [45–47]. By integrating deep learning, DRL significantly expands the applicability of RL to real-world problems where traditional methods would be computationally infeasible.

**2.3 Co-Design of Robots With Reinforcement Learning.** Given DRL's success in learning complex control policies, it has been increasingly adopted in robotic co-design, where the goal is to jointly optimize both the physical morphology of robots and their control policies [23–26,48,49]. This integrated perspective enables robots to evolve their mechanical structure and decision-making policies in tandem, yielding systems that can better adapt and perform in dynamic and uncertain environments [50].

A widely used strategy for solving CCD problems is to formulate them as a bi-level architecture. In this setup, the inner loop optimizes control policies for a fixed design candidate, while the outer loop updates the morphology through search methods such as evolutionary algorithms or Bayesian optimization [25,51–53]. This hierarchical decomposition allows the use of tailored solvers for each layer, accommodating discrete or nondifferentiable design spaces. However, bi-level methods are often computationally demanding and inefficient in sample usage, since every design modification requires retraining the controller from scratch [22].

Since DRL has shown great potential in optimizing control policies, its application to robot co-design, where both the physical hardware and control policies for the robots are optimized simultaneously, has gained increasing attention [23–26,48,49]. By extending the capabilities of DRL, robotic co-design focuses on developing algorithms that enable their mechanical structures to evolve and optimize along with their control strategies. Studies show that these robots have higher performance and adaptability in dynamic environments [50]. A new paradigm within this field is the concept of “hardware as policy,” where the robot's mechanical design is considered analogous to a control policy and optimized jointly with its computational counterpart using DRL [23,24]. By modeling hardware as auto-differentiable computational graphs, gradient-based algorithms from policy optimization can be used to efficiently cooptimize mechanical parameters alongside neural network weights and biases [25,26,48]. Such frameworks have been widely shown to scale to high-dimensional continuous control and design problems [25].

### 3 Methods

Traditional CCD methods rely on predefined mathematical models, limiting their adaptability to real-time variability in dynamic environments. This lack of adaptability leads to performance degradation in applications like vehicles and wind turbines due to wear and environmental changes [33]. To overcome this, the proposed CCD method uses DT frameworks for real-time data collection, continuous model updates, and informed decision-making on system design and policy updating. By integrating operational data, UQ, and DRL, our approach can iteratively refine the system model and improve designs for subsequent generations. This integration enhances the performance, adaptability, and robustness of these systems, which are difficult to achieve with traditional CCD methods.

It is important to note that, in stochastic and environment-driven settings, a unique ground-truth optimum is not well-defined unless explicit disturbance bounds or probability distributions are prescribed. In contrast to classical robust or stochastic CCD formulations that assume known uncertainty structures [9–11], the present work focuses on system adaptation under complex, potentially nonstationary disturbances that are learned from data. Consequently, performance is evaluated under consistent scenario realizations rather than against a single deterministic reference solution.

**3.1 Overview of Proposed Framework.** This article introduces a CCD framework for DTs that leverages DRL within a multi-generation design paradigm. The proposed framework enables the DT to evolve iteratively through successive learning cycles by continuously integrating data collected from physical systems during operation. Across generations, operational data is used to refine the digital model, incorporating quantile-based learning to account for

environmental uncertainties and policy learning to improve decision-making in dynamic and uncertain conditions. Prior to advancing to the next generation, the updated digital model, which is better aligned with real-world behaviors and variability, is employed to re-optimize the CCD problem. This iterative refinement ensures that both system designs and control strategies become progressively more adaptive and efficient over time.

Figure 3 illustrates the conceptual workflow of the proposed framework, composed of four steps. Prior to *Step 1*, an initialization stage is performed to construct the initial policy  $\pi$  and value function  $V_0$  using data generated by optimal control theory and Latin hypercube sampling (LHS), which provides a crucial foundation for guiding early learning and improving sample efficiency, shown in Fig. 3(b). In *Step 1*, shown in Fig. 3(c), the first CCD optimization is performed using the virtual model  $\mathcal{M}_0$ , initial physical design  $\mathbf{p}_0$ , and initial control policy  $\pi_0$ . Solving this problem yields the optimal system design  $\mathbf{p}_1$  and control policy  $\pi_1$ , which are deployed in the real environment in *Step 2* (referred to as Generation 1), shown in Fig. 3(d). During this phase, physical data are collected, and the digital model and the DRL-based control policy are updated in real time, resulting in an updated model  $\mathcal{M}_1$  and an updated controller  $\pi_2$ , respectively.

In *Step 3*, shown in Fig. 3(e), the refined model  $\mathcal{M}_1$  is used to resolve the CCD optimization problem, generating a new system design  $\mathbf{p}_2$  and updated controller  $\pi_2$ . This is followed by *Step 4*, where the new design and policy are implemented in the physical environment, initiating another cycle of data collection and model refinement, leading to  $\mathcal{M}_2$ . This process continues iteratively across generations, enabling continuous adaptation and performance improvement of both the system and its controller.

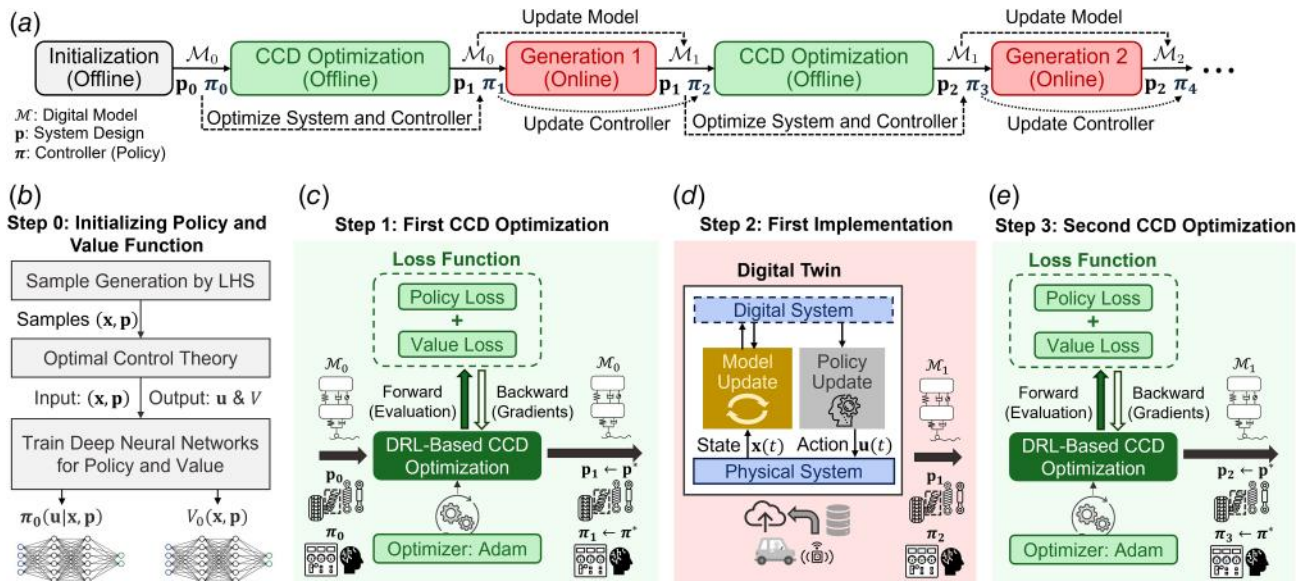
The proposed framework is particularly justified in systems where uncertainty is nonstationary and deployment dependent, such as evolving environmental conditions, user-specific operating patterns, or lifecycle degradation. While high-fidelity models can reduce structural errors under assumed scenarios, they are typically constructed offline and do not adapt to environment-specific variations encountered during operation. In contrast, the DT-enabled multigeneration approach enables continual data assimilation and model refinement, making it well suited for dynamic systems whose behavior evolves over time.

We note that the term “generation” in Fig. 3 does not imply repeated full-scale physical prototyping at every iteration. In many practical settings, model updating and controller refinement occur digitally using operational data collected after deployment. Physical redesign is performed selectively and strategically, for example when releasing a new product version or implementing major upgrades. Thus, the proposed framework aligns with lifecycle-based product development, where data from deployed systems informs systematic improvements across product generations rather than requiring continuous hardware reconstruction.

**3.2 Illustrative Example.** To demonstrate the proposed framework and its individual steps, we consider a simplified illustrative example. The system dynamics are described by a discrete-time linear model:

$$\mathbf{x}_{k+1} = \underbrace{\begin{bmatrix} 0.8 & 0.5 \\ 0.5 & 0.6 \end{bmatrix}}_{\mathbf{A}} \mathbf{x}_k + \underbrace{\begin{bmatrix} 0.5 \\ p \end{bmatrix}}_{\mathbf{B}} u_k + \mathbf{w}_k \quad (5)$$

where  $\mathbf{x}_k$  denotes the state vector at time-step  $k$  constrained by  $\mathcal{X} := \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 \in [-10, 5], x_2 \in [-5, 2]\}$ , and  $u_k \in \mathcal{U} := \{u \in \mathbb{R} \mid -1 \leq u \leq 1\}$  is the bounded control input. In implementation, sampled control actions are clipped to this admissible interval to ensure that input constraints are strictly satisfied. The term  $\mathbf{w}_k \in \mathbb{R}^2$  represents process noise or disturbances at time  $k$ , and the scalar parameter  $p \in [0.5, 2.0]$  serves as a tunable physical design variable within the input matrix  $\mathbf{B}$ .



**Fig. 3 Proposed multigeneration CCD framework for DT systems. (a) Overview of the CCD process across generations. (b) Offline initialization using LHS and optimal control theory to train initial policy  $\pi_0$  and value function  $V_0$ . (c) First CCD optimization using gradient-based DRL with auto-differentiation. (d) Online deployment and data collection in Generation 1, where real-time feedback is used to update the digital model and policy. (e) Second CCD optimization using updated models to improve performance in subsequent generations.**

The parameter  $p$  plays a critical role in the system's controllability. While increasing  $p$  enhances the system's responsiveness to control inputs, it may also amplify the sensitivity to control-induced errors, thus introducing a trade-off in the design process. The initial plant parameter was set to  $p_0 = 1.0$ , a nominal value commonly adopted for this benchmark system from Ref. [54].

The objective is to regulate the system to the origin. To this end, we define the reward function as

$$r_{k+1} = -\mathbf{x}_{k+1}^\top \mathbf{Q} \mathbf{x}_{k+1} - 0.1 u_k^2 \quad (6)$$

where  $\mathbf{Q} = \mathbf{I}_{2 \times 2}$  is the identity matrix used to penalize deviations from the origin in the state space, and the second term penalizes control effort.

Prior to system deployment, no physical data or prior environmental knowledge is assumed. To simulate real-world uncertainty, the disturbance  $\mathbf{w}_k$  is modeled as a zero-mean multivariate Gaussian noise process, given by

$$\mathbf{w}_k \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0.1^2 & 0 \\ 0 & 0.2^2 \end{bmatrix}\right), \quad \forall k \in \mathbb{N}_{\geq 0} \quad (7)$$

It is noted that the system in Eq. (5) is inherently unstable since  $\|\mathbf{A}^k\| \rightarrow \infty$  as  $k \rightarrow \infty$ . This instability highlights the necessity of an effective control strategy in conjunction with a well-chosen physical design  $p$  to stabilize the system, emphasizing the importance of proper policy initialization prior to learning.

**3.3 Step 0: Initialization.** A key challenge in DRL-based CCD lies in the initialization of the policy and value function approximators, typically represented by deep neural networks. Traditional methods often require solving multiple optimal control problems for diverse initial conditions, which is computationally demanding and may not scale well to high-dimensional design spaces. Poor initialization can also lead to unstable training and slower convergence, which reduces the practicality of DRL in real-time engineering applications.

To address this challenge, we adopt a sample generation strategy inspired by the set invariance and explicit MPC framework proposed by Chen et al. [55]. Their approach enables efficient sampling of feasible state-action pairs that satisfy system constraints

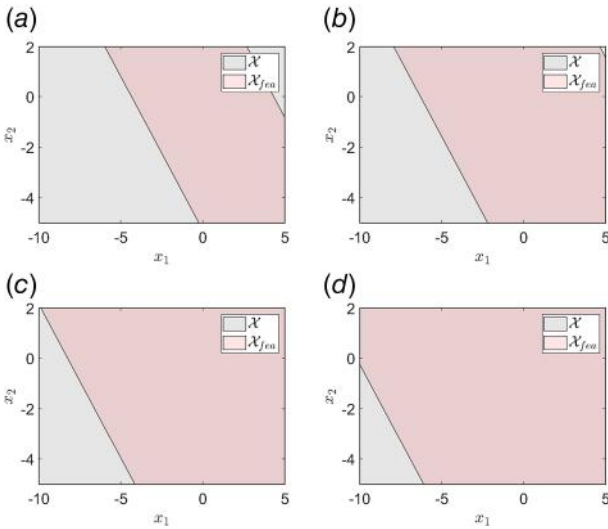
and approximate optimality, without the need to repeatedly solve optimization problems. This provides a data-driven yet control-theoretically grounded way to initialize neural network-based controllers.

**3.3.1 Incorporating Physical Design Parameters.** Inspired by Refs. [23,25,26,56], the input space extends beyond system states to include physical design variables  $\mathbf{p}$  in our CCD setting. To accommodate this, we generalize the method of Chen et al. [55] by introducing physical parameters  $\mathbf{p}$  directly into the input of the policy and value functions, yielding  $\pi(\mathbf{u}|\mathbf{x}, \mathbf{p})$  and  $V(\mathbf{x}, \mathbf{p})$ . For the illustrative example, since the physical parameter is a scalar  $p$ , we only add one more dimension to the sampling space. We will generate triplets  $(\mathbf{x}, p, u)$  where  $u$  is obtained by following Sec. 3.3.2. These samples provide a structured basis for initializing the neural networks, improving stability and constraint satisfaction in early training phases.

**3.3.2 Feasible Sample Generation via Invariant Set Computation.** To ensure that all generated samples are within regions of safe operation and satisfy both state and input constraints, the sample generation procedure is introduced and detailed in the following steps:

- (1) Define the discrete-time linear system  $\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k$ , where  $\mathbf{A}$ ,  $\mathbf{B}$ , or both depend on the physical parameters  $\mathbf{p}$ . For the illustrative example in Eq. (5), only  $\mathbf{B}$  depends on the physical parameter  $p$ .
- (2) Sample a set of candidate points  $(\mathbf{x}, \mathbf{p})$  uniformly within predefined bounds using LHS.
- (3) For each sample, compute its feasible  $\mathcal{X}_{\text{fea}}(\mathbf{p})$  using tools such as MPT3 [57]. Figure 4 visualizes  $\mathcal{X}_{\text{fea}}(p)$  for the illustrative example (Eq. (5)) with  $p = 0.5, 1.0, 1.5,$  and  $2.0$ .
- (4) Accept the sample  $(\mathbf{x}, \mathbf{p})$  only if  $\mathbf{x} \in \mathcal{X}_{\text{fea}}(\mathbf{p})$ .
- (5) For each accepted pair, assign a control input  $\mathbf{u}$  using the corresponding explicit MPC policy (detailed in Ref. [55]).

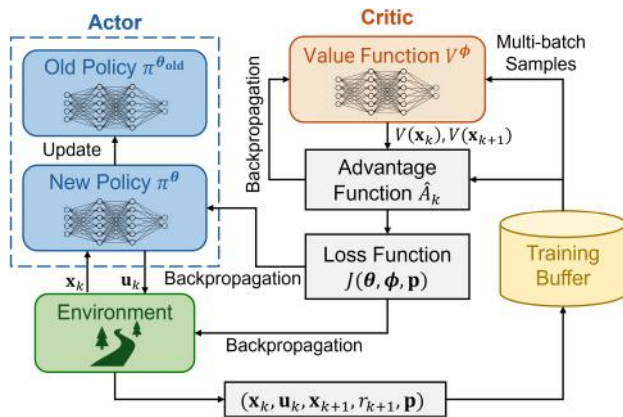
For the illustrative example (Eq. (5)), 300 samples  $(\mathbf{x}, p, u)$  are generated and then used to pretrain the DNNs that approximate the policy and value functions. Both networks are implemented as fully connected feedforward architectures with four hidden layers. Each network takes a 3-dimensional input vector: comprising the system



**Fig. 4 Visualization of  $\mathcal{X}_{\text{fea}}$  of the illustrative example for different values of  $p$ . The area  $\mathcal{X}_{\text{fea}}$  in each plot represents the state feasible region with the specified value of  $p$ , which expands as  $p$  increases. This trend reflects improved controllability of the system due to the larger entries in matrix  $\mathbf{B}$ , allowing a broader set of states to be driven to the origin under the state and control input constraints. (a)  $p = 0.5$ . (b)  $p = 1.0$ . (c)  $p = 1.5$ . (d)  $p = 2.0$ .**

state  $\mathbf{x} = [x_1, x_2]^\top$  and the scalar physical design parameter  $p$ . The hidden layers are configured with 32, 32, 16, and 16 neurons, respectively, with the first layer using a Tanh activation function to facilitate early nonlinear learning. The value network outputs a scalar estimate of the expected return, while the policy network is split into two branches: one for predicting the mean control action and the other for the standard deviation of the Gaussian policy. The mean network is pretrained using the sampled control actions from explicit MPC, while the standard deviation branch is initialized with small constant bias values (e.g., 0.01) and zero weights to allow for controlled exploration at the start of training.

This initialization technique significantly reduces the reliance on online optimal control solvers for pretraining and provides high-quality, constraint-satisfying data. By incorporating physical



**Fig. 5 Flowchart of the proposed DRL-based CCD optimization using the PPO algorithm. Unlike conventional PPO, this framework integrates the physical system parameter  $p$  as part of the input space and enables gradient-based updates of both the control policy and the physical design. This is achieved by treating the environment dynamics as differentiable and enabling backpropagation through the environment to the loss function.**

design variables directly into the input space, the method ensures that both structural and control-related decisions are grounded in feasible operation, enabling smoother and more sample-efficient DRL training.

**3.4 Step 1: Deep Reinforcement Learning-Based Control Co-Design Optimization.** Following the initialization of the neural network approximators, we present a DRL framework based on to solve the CCD problem, shown in Fig. 5. In this study, PPO [44] is adopted as the DRL optimizer due to its training stability and suitability for continuous control problems. Unlike value-based methods such as DQN [58], which are designed for discrete action spaces, the active control problems considered here involve continuous control inputs. PPO, as an actor-critic method with a clipped surrogate objective, provides stable policy updates while allowing joint optimization of policy parameters and physical design variables. The clipped objective prevents large policy shifts, which is particularly important in CCD settings where unstable updates could disrupt the coadaptation process between plant and controller. Other actor-critic methods could also be integrated into the proposed framework; however, the conceptual contributions of this work are independent of the specific policy gradient variant employed.

Unlike traditional PPO implementations, which focus solely on optimizing policy and value network parameters, our method simultaneously updates the physical design variables  $\mathbf{p}$  along with the policy parameters  $\theta$  and value function parameters  $\phi$  during training. This tight coupling allows both the physical system and the control policy to coadapt through reinforcement learning.

To achieve this, we formulate a joint optimization problem that minimizes a loss function composed of a clipped surrogate policy objective and a value prediction error:

$$\begin{aligned} \min_{\theta, \phi, \mathbf{p}} J(\theta, \phi, \mathbf{p}) &= \mathbb{E}_k \left[ \underbrace{-\min(\rho_k(\theta, \mathbf{p}) \hat{A}_k, L^{\text{CLIP}}(\rho_k(\theta, \mathbf{p}), 1 - \epsilon, 1 + \epsilon) \hat{A}_k)}_{\text{Policy Loss}} \right. \\ &\quad \left. + c_v \cdot \underbrace{L_{\text{SmoothL1}}(V^\phi(\mathbf{x}_k, \mathbf{p}), \hat{V}_k)}_{\text{Value Loss}} \right] \end{aligned} \quad (8)$$

where  $\hat{A}_k$  is the advantage function at time-step  $k$ , which estimates how much better an action is compared to the expected value of a state based on the reward function defined in Eq. (6), and  $L^{\text{CLIP}}(\rho_k(\theta, \mathbf{p}), 1 - \epsilon, 1 + \epsilon)$  is a function clipping  $\rho_k(\theta, \mathbf{p})$  within  $(1 - \epsilon, 1 + \epsilon)$ ,  $\rho_k$  denotes the probability ratio between the new and old policies:

$$\rho_k(\theta, \mathbf{p}) = \frac{\pi^\theta(\mathbf{u}_k | \mathbf{x}_k, \mathbf{p})}{\pi^{\theta_{\text{old}}}(\mathbf{u}_k | \mathbf{x}_k, \mathbf{p}_{\text{old}})} \quad (9)$$

$c_v$  is the coefficient for value loss, and we use Smooth L1 loss to define the value loss:

$$L_{\text{SmoothL1}}(a, b) = \begin{cases} \frac{1}{2}(a - b)^2, & \text{if } |x - y| < 1 \\ |a - b| - \frac{1}{2}, & \text{otherwise} \end{cases} \quad (10)$$

The inclusion of  $\mathbf{p}$  in both policy and value networks enables the optimizer to compute gradients not only with respect to the network weights but also with respect to the physical design parameters. To facilitate efficient gradient computation, we utilize automatic differentiation via the PyTorch autograd engine [59], which builds the computational graph and performs backpropagation for all trainable parameters, including  $\mathbf{p}$ . This allows for joint, end-to-end

updates during the learning process, enhancing sample efficiency and enabling exploration across both design and control spaces.

Figure 6 presents the training histories of the average returns and the system parameter  $p$  for the numerical example introduced in Sec. 3.2. The DRL-based CCD optimization was run for 10,000 epochs using a learning rate of  $10^{-5}$ . The remaining PPO hyperparameters are shown as follow: discount factor  $\gamma = 1$ , generalized advantage estimation (GAE) parameter  $\lambda = 0.9$ , clipping parameter  $\epsilon = 0.2$ , training epochs  $K = 10$ , rollout length 10, buffer size 10, and minibatch size 32. The training process required approximately 7 min on a machine with an AMD EPYC 7413 24-Core CPU (2.65 GHz), 1 TB RAM, and a 64-bit Linux operating system. As shown in the figure, the system parameter  $p$  evolves over the course of training, reflecting the coadaptation of the physical system and its control policy. Specifically, Fig. 6(b) illustrates how  $p$  initially increases to enhance system controllability. However, as excessively large values of  $p$  can amplify sensitivity to control inputs—especially if the policy is not yet fully trained—the learning algorithm occasionally reduces  $p$  to balance performance gains with robustness.

Figure 7 visualizes the mean of the learned policy  $\pi(u|\mathbf{x}, p_1)$  for the final optimized system parameter  $p_1 = 1.1513$ . Since the policy is stochastic by design, the control action may vary upon each execution. For visualization, we plot the mean of the distribution; in deployment, however, the policy continues to sample actions probabilistically.

To evaluate the performance of the optimized design, we simulate 1000 independent trajectories. The resulting state responses exhibit reduced settling times and smaller variations in overshoot compared to those before optimization, shown in Fig. 8. Quantitatively, the average return improves from  $-86.328$  (preoptimization) to  $-69.783$  (postoptimization), while the standard deviation decreases from 26.369 to 11.994. These results highlight that the proposed DRL-based CCD method successfully enhances dynamic performance by jointly optimizing the physical system parameter and control policy.

The current formulation relies on gradient-based updates and is therefore limited to continuous design variables, consistent with most CCD approaches in the literature. However, many practical

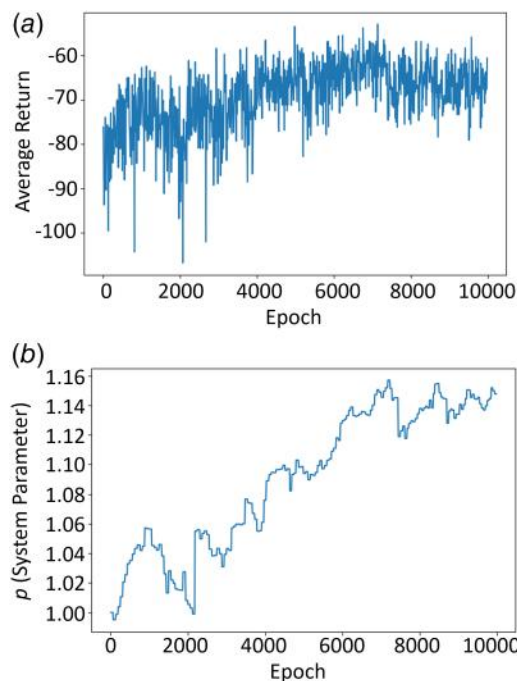


Fig. 6 Training history for numerical example with (a) average return and (b) system parameter  $p$

engineering problems involve discrete component specifications (e.g., standardized mechanical parts), and extending the framework to discrete or mixed-integer designs remains future work.

**3.5 Step 2: First Implementation (Generation 1).** It is noted that **Step 1** represents the conventional one-time CCD optimization performed before implementation and serves as the baseline design without operational data collected from environments for updating in subsequent generations. Prior to deploying the system in Generation 1, it is important to recognize the inherent mismatch between the digital model and the actual physical system. These discrepancies stem from unmodeled environmental variations (e.g., hysteresis and temperature-dependent damping in suspension components), environment-dependent disturbances (e.g., varying road conditions and driving behaviors), lifecycle effects such as wear and degradation, and other nonlinear effects that the nominal model does not fully capture.

To explicitly account for the mismatch between the digital and physical systems, we employ a discrepancy model trained on physical data collected during deployment. The concept of modeling structural discrepancy between computational models and reality has been extensively studied in the uncertainty quantification literature. Arendt et al. demonstrate that even calibrated high-fidelity models may exhibit systematic bias relative to observed data, and that separating model discrepancy from parameter uncertainty is essential for reliable prediction [60]. Motivated by this perspective, we incorporate a data-driven discrepancy function that captures residual bias between the nominal digital model and the physical system.

To reflect these modeling discrepancies, we construct a hypothetical physical system by augmenting the nominal model with additional disturbance and nonlinear terms:

$$\begin{aligned} \mathbf{x}_{k+1} = & \mathbf{A}\mathbf{x}_k + \mathbf{B}u_k + \mathbf{w}_k + \underbrace{\begin{bmatrix} 0.1 \\ -0.2 \end{bmatrix}}_{\text{Unknown Disturbances}} + U\left(\begin{bmatrix} -0.1 \\ -0.1 \end{bmatrix}, \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix}\right) \\ & + \underbrace{\left(\begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix} \text{diag}(\sin \mathbf{x}_k) + \begin{bmatrix} 0 & 0.1 \\ 0.1 & 0 \end{bmatrix} \text{diag}(\cos \mathbf{x}_k)\right)}_{\text{Unknown Nonlinearity}} \mathbf{x}_k \\ & + \underbrace{\left(\begin{bmatrix} 0.5 & 0 \\ 0 & 0 \end{bmatrix} \sin \mathbf{x}_k + \begin{bmatrix} 0 & 0 \\ 0 & 0.5p \end{bmatrix} \cos \mathbf{x}_k\right)}_{\text{Unknown Nonlinearity}} u_k \end{aligned} \quad (11)$$

where  $U(\text{lb}, \text{ub})$  represents a uniform distribution between the lower bound lb and upper bound ub. These additional terms simulate real-world uncertainties (such as additive disturbances and input/state-dependent nonlinearities) that are not incorporated in the original digital model defined in Eq. (5). This modeling gap emphasizes the importance of using a digital twin framework

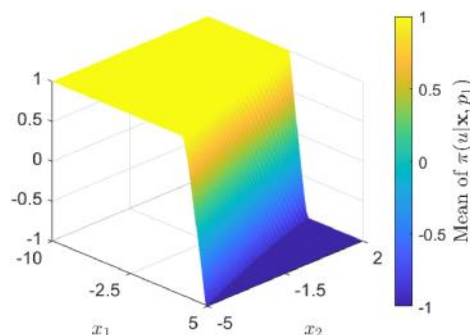
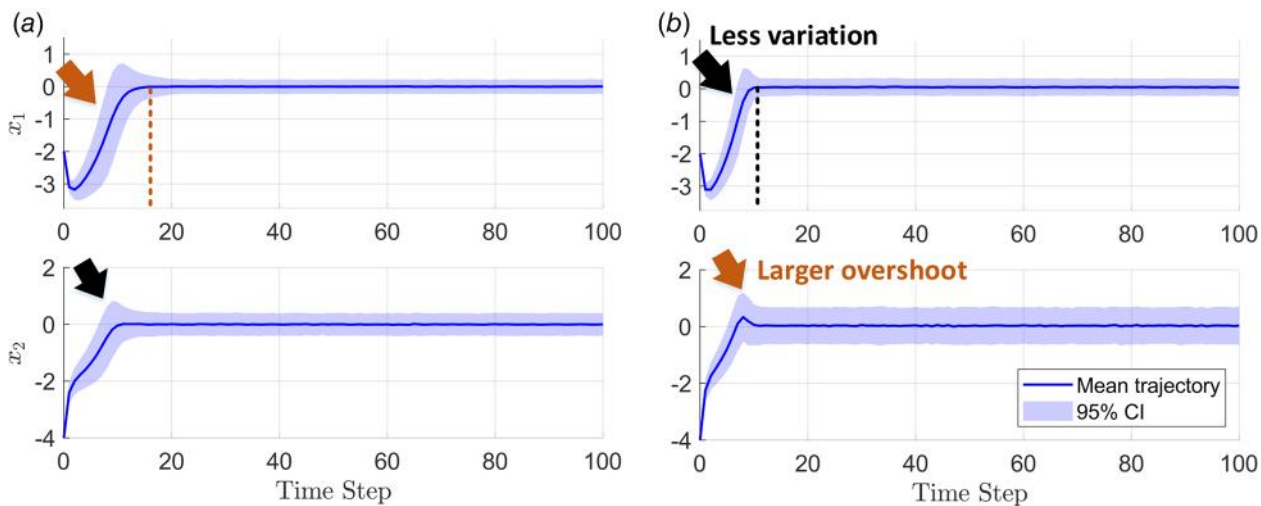


Fig. 7 Visualization of the mean of the optimal policy  $\pi(u|\mathbf{x}, p_1)$ , where  $p_1 = 1.1513$



**Fig. 8 State trajectories (a) before and (b) after the first CCD optimization for the numerical example with 1000 replicates. The shaded regions show the 95% confidence intervals.**

that can iteratively adapt based on physical observations collected during operation.

Despite the presence of discrepancies between the digital model and the real physical system, the DRL-based controller can maintain reasonable control performance when the mismatch between the nominal model and the physical system is moderate. This robustness stems from the inherently adaptive nature of learning algorithms, which are specifically designed to learn and evolve in dynamic and uncertain environments. During system deployment, the controller (also known as RL agent) continues interacting with the environment and collecting state–action trajectories from the physical system. These operational data serve two purposes: they are used to quantify model discrepancy and refine the digital twin model, and they can also be incorporated into the RL training process to further improve the control policy. By incorporating real-time data collected from the physical system, the RL agent continuously refines its policy, enabling it to make more informed and proactive decisions that maximize long-term rewards. However, as the discrepancy between the nominal model and the physical system grows, the performance of the originally optimized controller may degrade. This motivates the model updating and multigeneration CCD process proposed in this work. As a result, the controller can adapt to variations in system behavior, thereby mitigating the impact of modeling inaccuracies and unanticipated disturbances.

To explicitly account for the mismatch between the digital and physical systems, we employ a discrepancy model trained on physical data collected during deployment. This model captures deviations by leveraging UQ techniques. Specifically, quantile regression, a data-driven method for estimating conditional quantiles of a response variable. Quantile regression is well-suited for modeling uncertainties in dynamic systems, as it provides predictive intervals that characterize both systematic bias and aleatoric uncertainty (i.e., uncertainty inherent in the data) [61]. It also provides conformal estimation of the distribution which is free from assuming the distribution to be Gaussian, providing additional learning flexibility.

By selecting representative quantile levels, such as the 10th, 50th (median), and 90th percentiles, the discrepancy model captures the spread and central tendency of prediction errors. In the context of the numerical example, the model estimates the error at the next time-step as a function of the current state, control input, and current error:

$$\begin{bmatrix} \mathbf{e}_{k+1}^{\text{upper}} \\ \mathbf{e}_{k+1}^{\text{median}} \\ \mathbf{e}_{k+1}^{\text{lower}} \end{bmatrix}^T = \mathbf{f}_e(\mathbf{e}_k, \mathbf{x}_k, \mathbf{u}_k) \quad (12)$$

where  $\mathbf{e}_k := \mathbf{x}_k - \bar{\mathbf{x}}_k$  denotes the deviation between the actual state  $\mathbf{x}_k$  and the nominal predicted state  $\bar{\mathbf{x}}_k$  from the digital model, computed using  $\bar{\mathbf{x}}_{k+1} = \mathbf{A}\bar{\mathbf{x}}_k + \mathbf{B}\mathbf{u}_k$ . The predicted quantiles  $\mathbf{e}_{k+1}^{\text{upper}}$ ,  $\mathbf{e}_{k+1}^{\text{median}}$ , and  $\mathbf{e}_{k+1}^{\text{lower}}$  represent the range of likely deviations in the next state.

It is important to note that measurement noise and limited data introduce both aleatoric and epistemic uncertainty into the learning process [62]. The proposed framework does not assume that structural model discrepancy necessarily dominates these effects. Rather, it is most beneficial when residual model bias and environment-induced mismatch remain nonnegligible after accounting for nominal modeling and measurement noise. In this work, quantile regression captures predictive uncertainty, while multigeneration data assimilation progressively reduces epistemic uncertainty over time.

In this formulation, we assume that the primary sources of uncertainty stem from the external environment, rather than variability in the physical system parameter  $p$ . Therefore, the discrepancy model is designed to be independent of  $p$ . Figure 9 illustrates the predicted quantiles (10th, 50th, and 90th percentiles), capturing the uncertainty bounds of the model prediction. The model achieves a root mean square error of 0.1355 on the validation dataset, indicating reliable accuracy in capturing the discrepancy between the real and nominal systems.

The model correction step implicitly assumes that the discrepancy between the nominal digital model and the physical system is sufficiently structured and persistent to justify updating. If model discrepancies are negligible relative to measurement noise and learning uncertainty, multigeneration updating may yield limited improvement. The proposed framework is therefore particularly suited for systems operating in uncertain, evolving, or user-dependent environments, where structural mismatch accumulates over deployment cycles and cannot be fully mitigated by offline modeling alone.

**3.6 Step 3: Second Control Co-Design Optimization (Generation 2).** Instead of continuing with the nominal digital model subject to stochastic disturbances as defined in Eq. (5), we refine the model by incorporating the learned discrepancy function  $\mathbf{f}_e$  obtained in Step 2. The updated dynamics are given by

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{e}_{k+1}^{\text{median}} \quad (13)$$

where  $\mathbf{e}_{k+1}^{\text{median}}$  denotes the median predicted quantile of the error, as estimated by the discrepancy model in Eq. (12). To reflect this

model refinement in the learning process, the reward function is also updated as follows:

$$r_{k+1} = -\mathbf{x}_{k+1}^T \mathbf{Q} \mathbf{x}_{k+1} - 0.1u_k^2 - \Delta \mathbf{e}_{k+1}^T \mathbf{Q}_q \Delta \mathbf{e}_{k+1} \quad (14)$$

where  $\Delta \mathbf{e}_{k+1} := \mathbf{e}_{k+1}^{\text{upper}} - \mathbf{e}_{k+1}^{\text{lower}}$  captures the predicted uncertainty range (i.e., quantile width) at time-step  $k+1$ , and  $\mathbf{Q}_q = 0.1\mathbf{Q}$  is a weighting matrix applied to penalize uncertainty. This revised reward function encourages the RL agent not only to maximize performance but also to minimize predictive uncertainty, thereby promoting more robust design and control decisions.

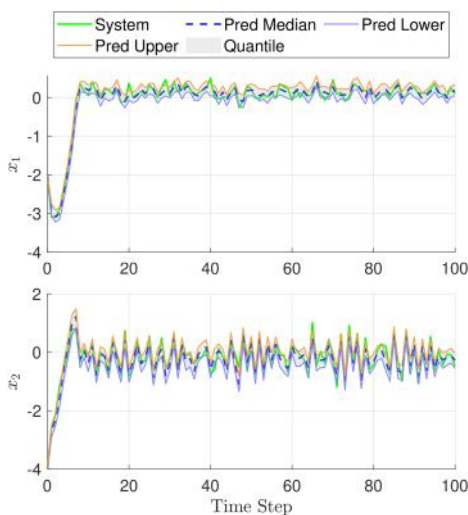
In *Step 3*, we perform a second CCD optimization using the updated digital model and reward formulation. The new optimized physical design parameter is  $p_2 = 1.2947$ . As previously discussed, increasing  $p$  improves controllability but requires a sufficiently capable controller. The upward shift in  $p$  reflects the coadaptation of the physical system and the RL policy, indicating that the controller has matured in its ability to manage a more sensitive system.

To evaluate the updated design, we simulate 1000 rollout trajectories. The average return improves from  $-113.817$  to  $-108.301$ , while the standard deviation decreases from 20.115 to 14.681 after the second optimization. Although the gain in expected return is modest, the significant reduction in variability demonstrates improved robustness of the closed-loop system.

The proposed framework follows a multigeneration design philosophy, where each generation involves cooptimization of system and control parameters informed by real-world data and model refinement. In this article, we illustrate the initialization stage and the first three steps (first optimization, first physical implementation, and second optimization); however, the framework is inherently extensible. The number of generations required is problem dependent and should be determined based on performance objectives and application constraints. Designers are encouraged to iteratively apply additional generations as needed to achieve desired levels of performance, reliability, and adaptability.

#### 4 Engineering Study: Active Suspension System

Unlike passive suspensions, active suspension systems enhance passenger comfort and vehicle stability by dynamically adjusting damping forces. However, their design requires a synergetic collaboration between the physical and control domains to improve comfort and stabilize the vehicle under various disturbances.



**Fig. 9 Visualization of the learned quantiles and the real system trajectories for the numerical example, where “time-step” represents  $k$  in Eq. (5)**

Although traditional proportional–integral–derivative controllers are simple to implement, they suffer from significant limitations, including a lack of adaptability and the need for precise tuning of control gains. While MPC, which optimizes real-time control actions by predicting future system responses, improves performance by accounting for future disturbances (shown in Chapter 8 of Ref. [7]), its applications are often limited by computational resources.

Due to the need for rapid decision-making for active suspension systems, such systems are chosen as the case study for demonstrating the proposed framework with DRL controllers which can provide near-instantaneous adjustments to suspension dynamics by evaluating the explicit policy. Additionally, once deployed in a real-world setting, the system generates a large volume of operational data, which can be leveraged to continuously train and update the controller, improving performance over time. Moreover, suspension systems operate in highly uncertain and dynamic environments, where road conditions, driving speed, and external disturbances vary unpredictably. DRL is well-suited for handling such uncertainties by learning adaptive control strategies that generalize across diverse conditions.

The quarter-car model of a vehicle suspension is shown in Fig. 10, where  $m_s$  and  $m_{us}$  are the sprung and unsprung masses, respectively,  $z_s$  and  $z_{us}$  are the vertical positions of the sprung and unsprung masses, respectively,  $z_0$  is the elevation of the road which is the excitation source to the system,  $k_t$  and  $c_t$  are tire stiffness and damping constant, respectively, and  $k_s$  and  $c_s$  are the coefficients of the spring and the damper (design variables for the physical system). The values of the parameters of the suspension system can be found in Ref. [4]. The dynamic equation is defined by

$$\begin{aligned} \underbrace{\begin{bmatrix} \dot{z}_{us}(t) - \dot{z}_0(t) \\ \ddot{z}_{us}(t) \\ \dot{z}_s(t) - \dot{z}_{us}(t) \\ \ddot{z}_s(t) \end{bmatrix}}_{\mathbf{\dot{x}}(t)} &= \underbrace{\begin{bmatrix} 0 & 1 & 0 & 0 \\ -\frac{k_{us}}{m_{us}} & -\frac{c_s}{m_{us}} & \frac{k_{us}}{m_{us}} & \frac{c_s}{m_{us}} \\ 0 & -1 & 0 & 1 \\ 0 & \frac{c_s}{m_s} & -\frac{k_s}{m_s} & -\frac{c_s}{m_s} \end{bmatrix}}_{\mathbf{A}_c} \underbrace{\begin{bmatrix} z_{us}(t) - z_0(t) \\ \dot{z}_{us}(t) \\ z_s(t) - z_{us}(t) \\ \dot{z}_s(t) \end{bmatrix}}_{\mathbf{x}(t)} \\ &+ \underbrace{\begin{bmatrix} 0 \\ -\frac{1}{m_{us}} \\ 0 \\ \frac{1}{m_s} \end{bmatrix}}_{\mathbf{B}_c} u(t) + \underbrace{\begin{bmatrix} -1 \\ \frac{c_t}{m_{us}} \\ 0 \\ 0 \end{bmatrix}}_{\mathbf{E}_c} \dot{z}_0(t) \end{aligned} \quad (15)$$

By selecting the sampling time  $T = 0.05$  s, the continuous-time dynamic model in Eq. (19) can be converted in a discrete-time version:

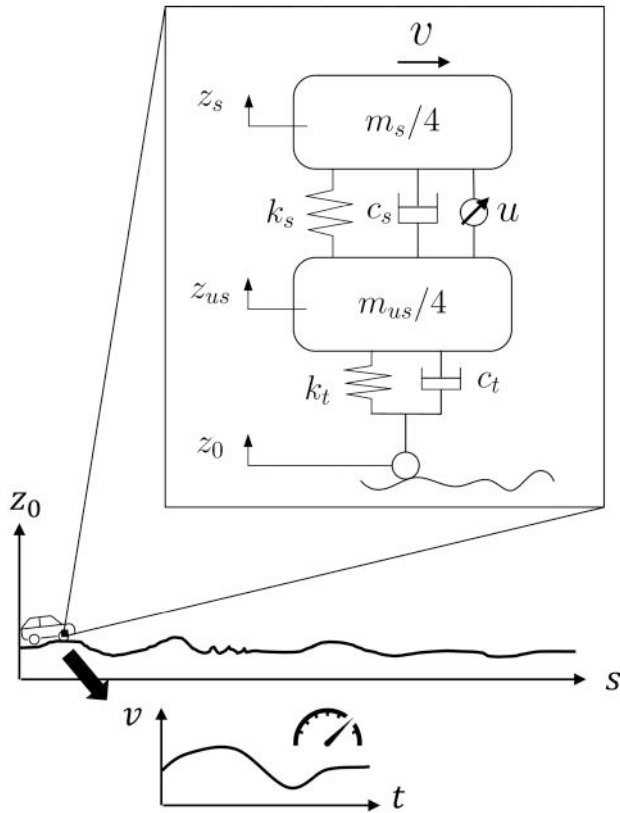
$$\mathbf{x}_{k+1} = \mathbf{A} \mathbf{x}_k + \mathbf{B} u_k + \mathbf{w}_k \quad (16)$$

where  $\mathbf{A} = e^{\mathbf{A}_c T}$  and  $\mathbf{B} = \int_0^T e^{\mathbf{A}_c \tau} \mathbf{B}_c d\tau$ . Although matrix  $\mathbf{B}_c$  is independent on  $k_s$  and  $c_s$ , matrix  $\mathbf{B}$  can be affected by their values after discretization. The disturbance vector  $\mathbf{w}_k$  is computed by

$$\mathbf{w}_k = \mathbf{E}_d \dot{z}_0(t_k) \quad (17)$$

where  $\mathbf{E}_d = \int_0^T e^{\mathbf{A}_c \tau} \mathbf{E}_c d\tau$  and  $\dot{z}_0(t_k)$  is the first derivative of the elevation at time  $t = t_k$ , which can also be expressed as  $\dot{z}_0(t_k) = (dz_0(s)/ds)v(t_k)$ , where  $s$  denotes distance. It means that the disturbance  $\dot{z}_0(t_k)$  depends on the road roughness and the vehicle speed  $v$ .

The task of the suspension system is to minimize vibrations and enhance ride comfort by regulating the state variables toward desired equilibrium points, effectively reducing oscillations caused by road roughness and other disturbances. Another important consideration is minimizing the control effort in order to



**Fig. 10 Quarter-car vehicle suspension model, modified from Ref. [7]. Many uncertain factors affect the suspension system performance from the environment like road conditions and driving speed. Note that the coefficients for the spring and the damper,  $k_s$  and  $c_s$ , are system parameters to be designed.**

reduce energy consumption. The reward function is defined by

$$r_{k+1} = -\mathbf{x}_{k+1}^T \mathbf{Q} \mathbf{x}_{k+1} - 10^{-6} u_k^2 \quad (18)$$

where  $\mathbf{Q} = \text{diag}([10, 1, 50, 5])$ .

#### 4.1 Step 1: Initial Control Co-Design Optimization.

To optimize the system with disturbances without prior knowledge of the environment, we assume the disturbance vector follows a Gaussian distribution, i.e.,  $\dot{z}_0(t_k) \sim \mathcal{N}(0, 0.3^2)$ . The initial values of the physical components are  $k_s = 27692.0$  N/m and  $c_s = 1906.5$  N s/m. For this case, we adopt similar DNN structures as described in the numerical example: fully connected feedforward networks with five hidden layers (16, 32, 32, 16, and 1 neurons) and Tanh activations. The input is 6-dimensional ( $x_1, x_2, x_3, x_4, k_s, c_s$ ), and the output is a scalar of control action  $u$ . The policy and value networks share the same architecture, with the policy consisting of two DNNs for the mean and standard deviation. The mean policy network is preinitialized using 3,000 samples, while the standard deviation network is initialized with zero weights and bias values of 0.01. The pretraining step is computationally inexpensive. Generating 3000 optimal state–design samples required approximately 3 min, and training the policy and value networks required an additional 3 min.

We use 10,000 epochs for the CCD optimization with the learning rate of  $10^{-5}$ . The learning rate was chosen to balance stability and convergence speed, preventing instability while enabling effective policy refinement, and 10,000 epochs were empirically selected to ensure sufficient training iterations for convergence without excessive computational cost, as fewer epochs led to underfitting while significantly more did not provide substantial additional improvements. The remaining PPO hyperparameters

are shown as follow: discount factor  $\gamma = 1$ , GAE parameter  $\lambda = 0.9$ , clipping parameter  $\epsilon = 0.2$ , training epochs  $K = 20$ , rollout length 20, buffer size 10, and minibatch size 32. The training process took approximately 3 h on a system equipped with an AMD EPYC 7413 24-core CPU (2.65 GHz), 1 TB of RAM, running a 64-bit Linux operating system. The values have been updated to  $k_s = 31,863.8$  N/m and  $c_s = 2067.1$  N s/m. The average returns before and after optimization are  $-87.429$  and  $-82.106$ , respectively, while the standard deviations of these 1000 trajectories before and after optimization are 11.932 and 12.059. The improvement looks trivial at this stage because the system behavior is simple with a predefined form of disturbances and the initial design is close to optimum so far. In the next step as the system is placed in the real environment and the physical data is continuously collected, we can get more information about the environmental conditions and capture some underlying physics.

**4.2 Step 2: First Implementation (Generation 1).** There are inherent discrepancies between the real system and its digital model, such as unknown environmental variations, and nonlinearities that our digital model does not fully capture. In the context of suspension systems, environmental disturbances come from unknown road conditions and driving behaviors. A road profile of several jumps, bumps, dents, and ramps with some noise was created over distances, while a speed profile with 15,000 steps was generated as a function of time.<sup>2</sup> The system nonlinearities exist in the mechanical components of the spring and damper. The dynamics of the real suspension system are defined as

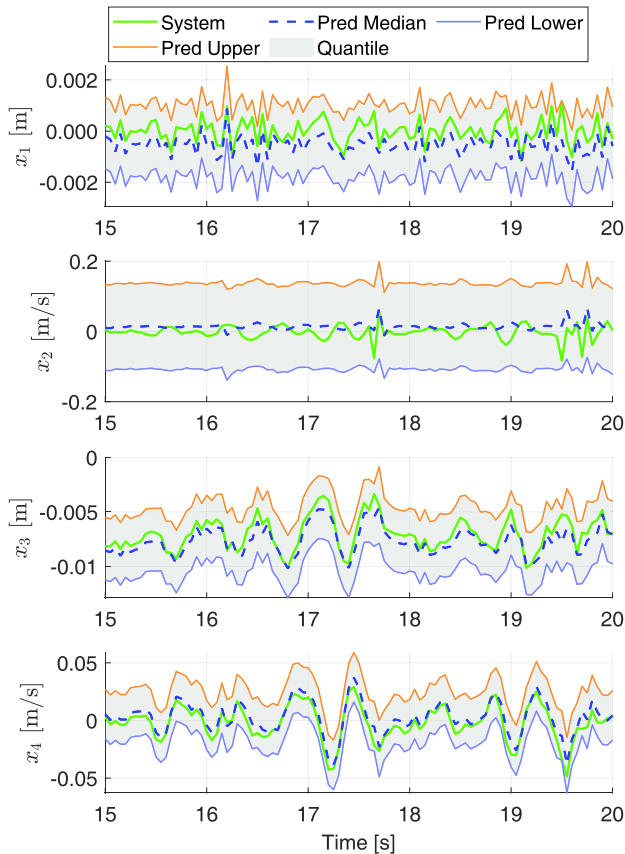
$$\dot{\mathbf{x}}(t) = \mathbf{A}_c \mathbf{x}(t) + \mathbf{B}_c u(t) + \mathbf{E}_c \dot{z}_0(t) + \underbrace{\begin{bmatrix} 0 \\ -\frac{1}{m_s} (k_{nl} x_3^3(t) + c_{nl} |x_4(t) - x_2(t)| (x_4(t) - x_2(t))) \\ 0 \\ \frac{1}{m_s} (k_{nl} x_3^3(t) + c_{nl} |x_4(t) - x_2(t)| (x_4(t) - x_2(t))) \end{bmatrix}}_{\text{Nonlinear terms of spring and damper}} \quad (19)$$

where  $k_{nl} = 0.01 k_s$  and  $c_{nl} = 0.01 c_s$  represent the nonlinear coefficients of the spring and damper, whose values follow a ratio of the coefficients  $k_s$  and  $c_s$ , and  $\dot{z}_0(t)$  is provided given the profiles of road elevations and driving speeds.

In practice, we can use sensors to measure the vertical movements and excitations from the road as well as the driving speeds in real time. The collected data can be used for updating the RL policy and the digital model. A discrepancy model, defined in Eq. (12), is trained using quantile regression. Figure 11 shows the learned quantiles with the real system trajectories, where the lower, median, and upper quantiles represent 10th, 50th, and 90th percentiles, respectively. The root mean square error for the validation dataset is 0.0343. The discrepancy model captures the uncertainty for most of the state variables with effectively learned quantiles. For the second state variable (tire vertical velocity), the predicted uncertainty band is relatively wider. This behavior is attributed to occasional abrupt velocity transitions induced by rough road excitations, which introduce locally high conditional variability and heteroscedastic effects. To remain robust in these regions, the quantile regression model produces wider prediction intervals, resulting in a more conservative uncertainty estimate.

**4.3 Step 3: Second Control Co-Design Optimization (Generation 2).** The main difference between the first and second CCD optimizations is that the second CCD optimization uses the updated dynamic model with the learned quantile,

<sup>2</sup>The codes and data for the road and speed profiles are provided [https://github.com/TsaiYK/ControlCoDesign\\_DigitalTwin\\_RL.git](https://github.com/TsaiYK/ControlCoDesign_DigitalTwin_RL.git)

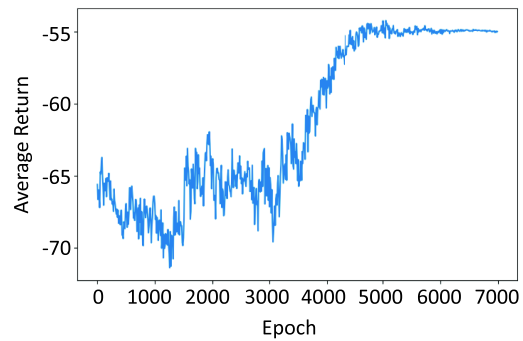


**Fig. 11 Visualization of the learned quantiles and the real system trajectories for the active suspension system**

shown in Eq. (13), and the updated reward function, shown in Eq. (14). With the physical data collected from Generation 1, the updated model more accurately captures environmental uncertainty through quantile learning, and there is no need for explicit assumptions about the disturbances. In this study, the weighting coefficient for the uncertainty penalty is set as  $\mathbf{Q}_q = 0.1\mathbf{Q}$ , consistent with the setting used in Eq. (14). This value provides a moderate penalty on predictive uncertainty while maintaining the primary emphasis on the nominal performance objective.

The return history from the RL process exhibits an overall upward trend shown in Fig. 12, indicating that the RL policy is progressively improving to maximize returns. The fluctuations observed in the trajectory are expected, as the decision-making process in RL is inherently stochastic. This variability arises from the trade-off between exploration and exploitation, which is crucial for operating in highly uncertain environments. For the active suspension system, the learning trend demonstrates a steady improvement in policy performance. Furthermore, it is important to note that the results demonstrate a clear upward trend and indicate consistent performance improvement and convergence of the policy over time. The steady increase and eventual stabilization of the average return suggest that the proposed method effectively learns a near-optimal co-design solution. This behavior supports the reliability of the learning process and highlights the robustness of the approach in handling the simultaneous optimization of control and physical design variables under uncertainty. The framework jointly optimizes the control strategy while adapting the physical parameters of the suspension system.

To provide a reference solution, a linear MPC framework is implemented based on the discretized state-space model of the quarter-car system. At each time-step, the control input is obtained by solving a finite-horizon optimization problem with horizon length  $N = 15$ , minimizing a quadratic cost composed of



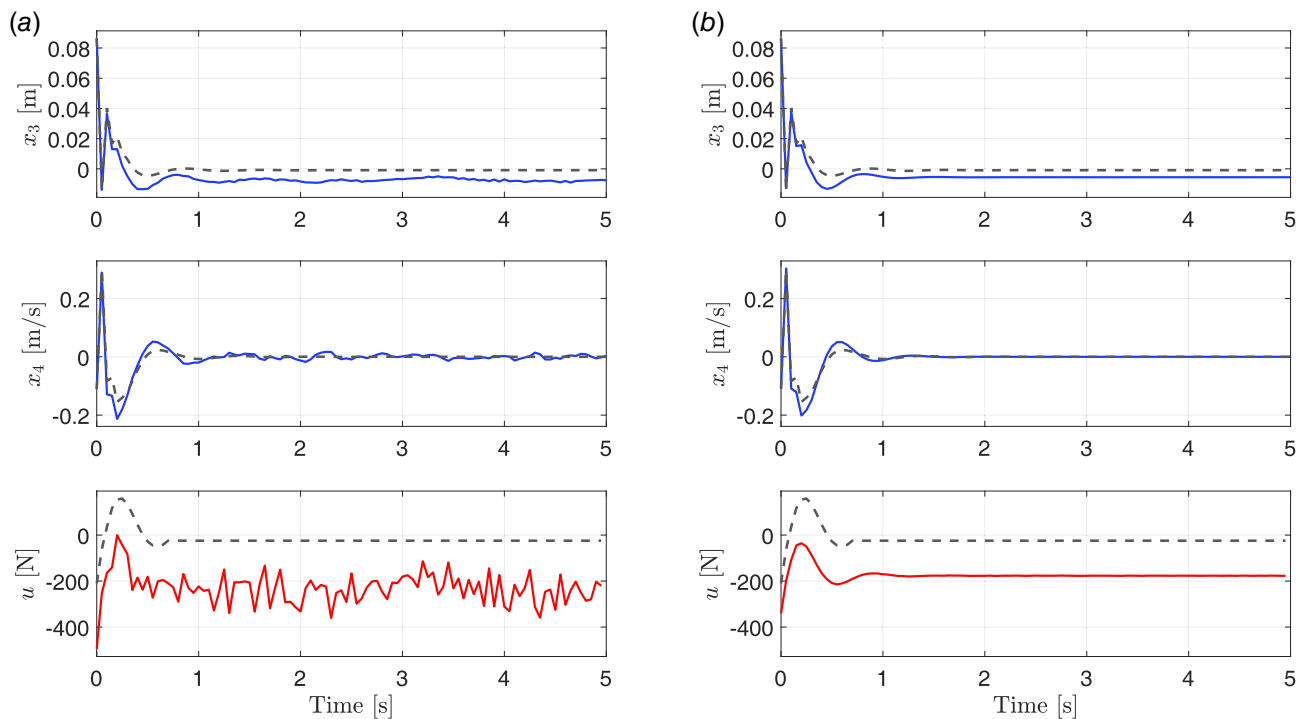
**Fig. 12 Training history for the active suspension system. The optimized values of the physical system parameters (the coefficients of the spring and damper) are  $k_s = 32290.6 \text{ N/m}$  and  $c_s = 2072.6 \text{ N s/m}$ .**

state penalties and control effort, consistent with the weighting matrices used in the RL formulation. The optimization is solved using a quadratic programming solver with input constraints  $u \in [-1000, 1000] \text{ N}$  and terminal state constraints defined by the admissible state bounds. The same system parameters ( $k_s$ ,  $c_s$ ) and initial conditions are used for all cases, and the MPC controller is evaluated in closed-loop using the nonlinear system dynamics to ensure a fair comparison with the proposed method.

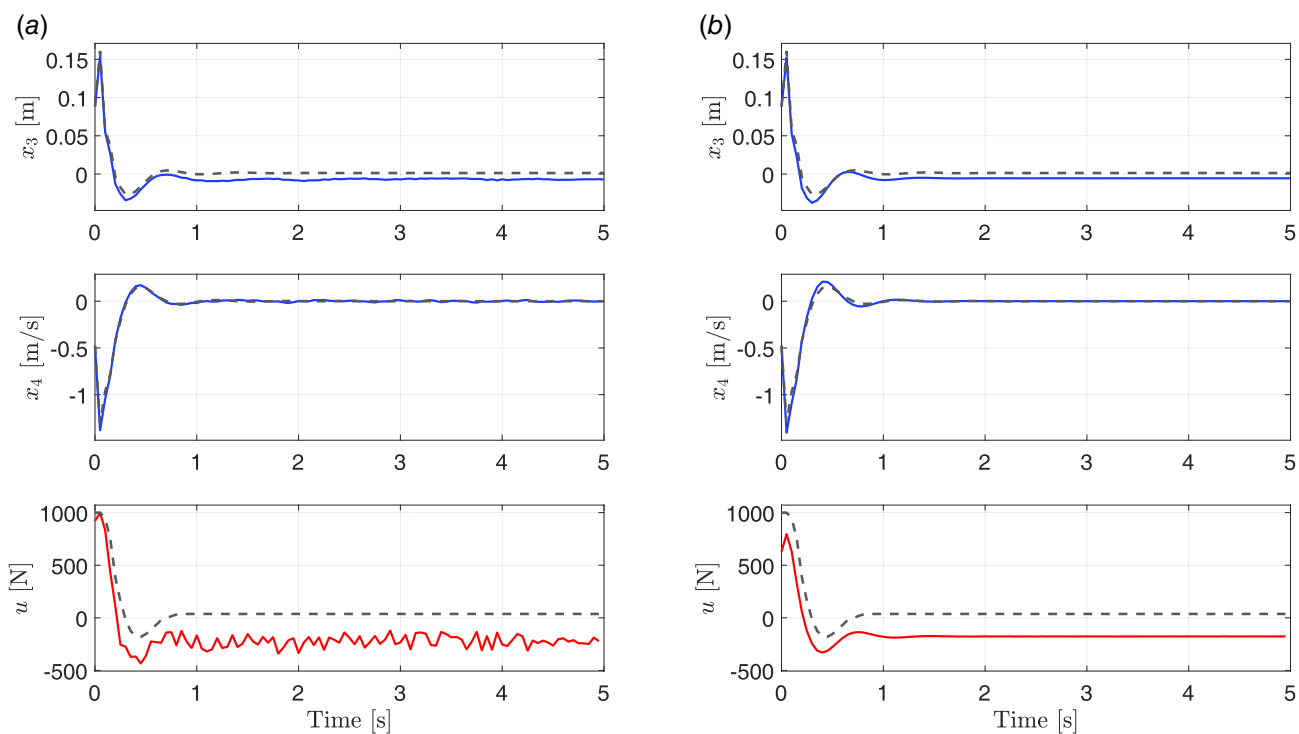
For the sake of simplicity, we refer to the solutions before and after the second CCD optimization as Generation 1 design (Gen-1) and Generation 2 design (Gen-2), respectively. The “before 2nd CCD” case represents a single-generation CCD baseline, where plant and controller are optimized using the nominal digital model without incorporating deployment data for model correction. The subsequent “after 2nd CCD” results isolate the benefit of the proposed multigeneration digital-twin update loop. To compare their performance, 200 samples of different initial conditions, we generate 200 samples of different initial conditions within the range  $\mathbf{x}_{ib} = [-0.5, -2, -0.2, -1]^T$  to  $\mathbf{x}_{ub} = [0.5, 2, 0.2, 1]^T$ . Both designs are simulated for 100 steps (equivalent to 5 s), and their dynamic performance is evaluated by computing the returns for all 200 trajectories. The mean return for the Gen-1 design is  $-26.3327$ , while the Gen-2 design achieves an improved mean return of  $-24.4252$ . The corresponding standard deviations are 19.6797 for Gen-1 and 19.7557 for Gen-2. These results indicate that the Gen-2 design achieves improved performance while maintaining a comparable level of robustness. This demonstrates that incorporating physical data and updating the digital model can enhance the design of autonomous systems.

To qualitatively assess the improvements achieved through the second CCD optimization, the dynamic responses of the Gen-1 and Gen-2 designs (corresponding to the system before and after optimization, respectively) are compared in Figs. 13–15. The responses are evaluated under different initial conditions. In addition, the results from a linear MPC controller are included as a reference. The linear MPC is constructed based on the nominal linear model and quadratic cost function, which can be efficiently solved to obtain a reference optimal solution under known model assumptions. The sharp increase in control input observed at the start of the trajectories ( $u$  in Figs. 13–15) is due to the use of nonequilibrium initial conditions. This setup was intentionally chosen to evaluate and compare the dynamic response and stabilization behavior of different designs. In real-world scenarios, systems typically begin near equilibrium, and such abrupt control actions are less likely to occur.

In Figs. 13–15, the Gen-2 design exhibits significantly more stable responses with reduced control effort. Compared to the linear MPC reference, Gen-2, the solution after updating the model and codesigning the physical system and controller with



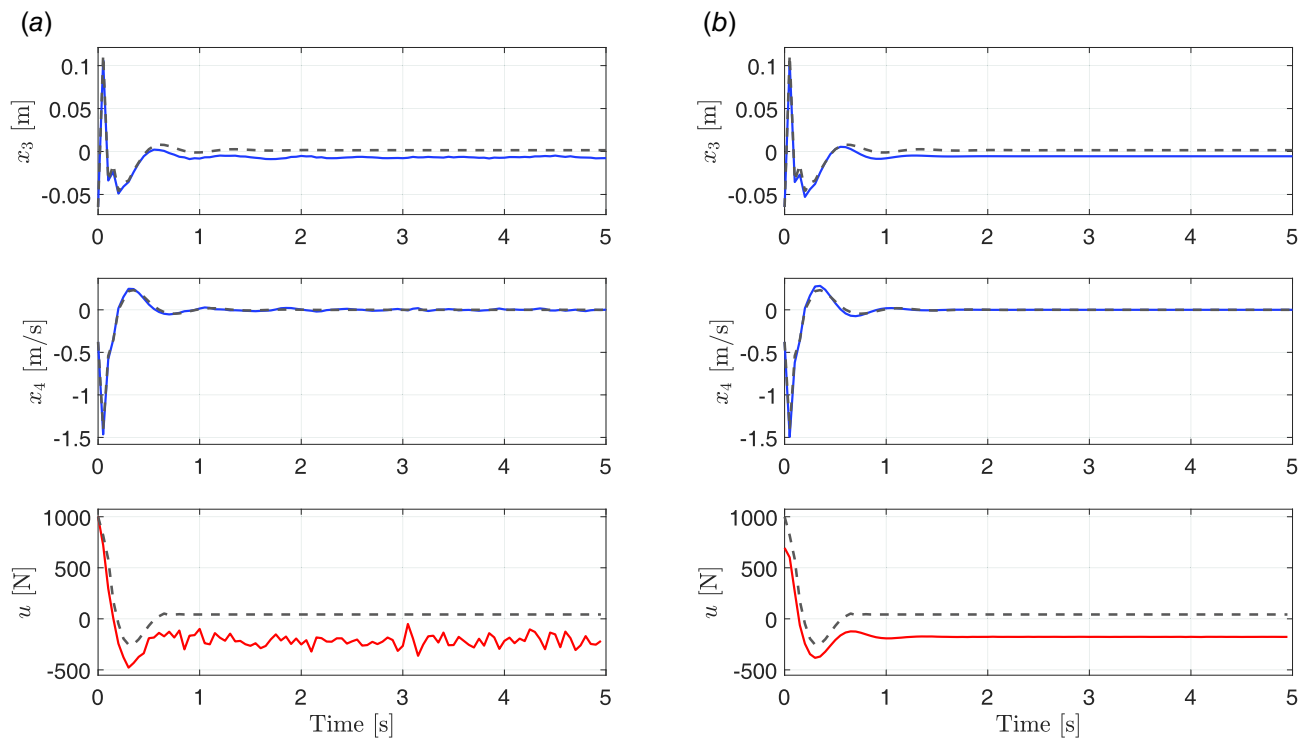
**Fig. 13** Third and fourth states and control trajectories of Gen-1 and Gen-2 designs, with the results by linear MPC as a reference (dashed gray curves), for the initial condition  $x = [-0.067, -0.178, 0.086, -0.110]^T$ . (a) Gen-1 and (b) Gen-2.



**Fig. 14** Third and fourth states and control trajectories of Gen-1 and Gen-2 designs, with the results by linear MPC as a reference (dashed gray curves), for the initial condition  $x = [0.086, -0.191, 0.088, -0.477]^T$ . (a) Gen-1 and (b) Gen-2.

collected data from Generation 1, achieves comparable stabilization performance, while maintaining smoother control trajectories in the presence of model discrepancies. It is important to note that, although the MPC solution provides a strong reference due to the linear dynamics and quadratic cost formulation, the real

system described in Eq. (19) includes nonlinear terms, and therefore the MPC result cannot be regarded as the true ground truth but rather as a high-quality benchmark. Moreover, the control trajectories in Gen-2 are noticeably smoother, indicating improved efficiency. It is also observed that the MPC inputs converge to



**Fig. 15** Third and fourth states and control trajectories of Gen-1 and Gen-2 designs, with the results by linear MPC as a reference (dashed gray curves), for the initial condition  $\mathbf{x} = [0.146, -0.208, -0.065, -0.377]^T$ . (a) Gen-1 and (b) Gen-2.

**Table 1** Returns for Generation 1 design (Gen-1) and Generation 2 design (Gen-2), with the results by linear MPC as a reference

| Cond | Linear MPC | Gen-1    | Gen-2           |
|------|------------|----------|-----------------|
| 1    | -1.4772    | -6.5324  | <b>-4.6694</b>  |
| 2    | -23.7094   | -29.4001 | <b>-27.6854</b> |
| 3    | -17.2486   | -22.4732 | <b>-20.9795</b> |

Note: Higher values indicate better performance, and bold entries demonstrate the superiority of Gen-2 over Gen-1.

zero (or near zero) in the steady-state, whereas the RL-based controllers maintain small but nonzero control actions. This behavior arises because the weight associated with control effort in the reward function is relatively small ( $10^{-6}$ ), leading the RL agent to prioritize state regulation over minimizing control usage. Furthermore, since the RL-based CCD framework jointly optimizes both control actions and physical design parameters, improvements in system performance are achieved not only through control but also through adjustments of the physical system, which may reduce the need for strictly minimizing control effort. This trend is further supported by the quantified results in Table 1, where the Gen-2 design consistently achieves improved returns compared to Gen-1 across all initial conditions, while remaining reasonably close to the linear MPC reference. In addition, while MPC requires solving an optimization problem online at each time-step, which may introduce computational latency in real-time applications, the computational cost of solving the MPC problem is not reflected in the trajectory results shown here. In contrast, the learned RL controller can generate control actions with negligible computational cost, making it more suitable for real-time deployment.

To further quantify the performance, the standard deviations of the steady-state responses,  $\sigma_{ss}$ , under three different initial conditions are computed and displayed in Table 2. Compared to Gen-1, Gen-2 achieves consistently lower values across both state variables and control actions. Notably, the reduction in the

**Table 2** Standard deviations of the steady-state trajectories  $\sigma_{ss}$  for Gen-1 and Gen-2 under three different initial conditions, with the results by linear MPC as a reference

| Cond | Metric      | Linear MPC | Gen-1   | Gen-2         |
|------|-------------|------------|---------|---------------|
| 1    | $x_3$ (m)   | 0.0001     | 0.0010  | <b>0.0001</b> |
|      | $x_4$ (m/s) | 0.0016     | 0.0078  | <b>0.0018</b> |
|      | $u$ (N)     | 0.0000     | 58.3250 | <b>1.1481</b> |
| 2    | $x_3$ (m)   | 0.0004     | 0.0014  | <b>0.0008</b> |
|      | $x_4$ (m/s) | 0.0028     | 0.0142  | <b>0.0064</b> |
|      | $u$ (N)     | 0.0000     | 53.8368 | <b>2.3714</b> |
| 3    | $x_3$ (m)   | 0.0005     | 0.0013  | 0.0017        |
|      | $x_4$ (m/s) | 0.0043     | 0.0151  | <b>0.0139</b> |
|      | $u$ (N)     | 0.0000     | 54.7468 | <b>2.8782</b> |

Note: Lower values indicate better performance, and bold entries demonstrate the superiority of Gen-2 over Gen-1.

fourth state response,  $x_4$ , is especially important, as this variable directly relates to passenger comfort by mitigating oscillatory motion. At the same time, the control effort  $u$  is substantially reduced, leading to smoother and more energy-efficient actuation. These results confirm that the Gen-2 design not only enhances comfort and robustness but also achieves significant energy savings in actuation.

## 5 Conclusion and Future Work

This work presents a lifecycle-oriented CCD framework that integrates DTs and DRL to address the challenge of designing dynamic systems under unpredictable uncertainties. The proposed approach extends CCD beyond single-cycle optimization by enabling DTs to evolve through real-time sensing, model updating, and adaptive re-optimization across successive generations. DRL accelerates real-time decision-making and controller adaptation, while multigeneration design leverages operational data from each deployed system to refine DT models, improve uncertainty

quantification through quantile regression, and inform the redesign of both physical and control components for subsequent generations. Together, these elements ensure that system performance and robustness improve progressively over the lifecycle. The key contributions of this work are: (1) extending CCD into a multigeneration framework that explicitly incorporates the full product lifecycle, (2) leveraging DTs to enable continuous model evolution and data-driven decision-making, and (3) employing DRL to accelerate adaptive control and enhance real-time responsiveness.

The proposed framework is demonstrated on an active suspension system, where DRL and automatic differentiation are employed to cooptimize the system parameters for the mechanical components and control policies. Quantile regression is used to learn the environmental uncertainty from physical data on road conditions and driving speeds after the deployment. The results indicate that the optimized solution after the second CCD optimization, which incorporates the updated model and control policy, achieves better dynamic responses, smoother control trajectories, and improved robustness compared to the initial design. A key advantage of this framework is its ability to cooptimize the physical system and its controller to ensure that performance and robustness improve over successive generations.

One limitation of the proposed approach is the need for many training epochs to ensure stable convergence when jointly optimizing the system and policy. This is due to the complexity of coadaptation and the goal of achieving generalizability. While the proposed DRL-based framework demonstrates stable convergence in the presented case studies, convergence is not guaranteed in general, particularly for high-dimensional and complex systems. Furthermore, the choice of reward function plays a critical role in guiding the learning process. Although the engineering objective is directly used as the reward in this work, more complex problems may require additional reward shaping to facilitate efficient learning and avoid suboptimal policies. In our recent work on full-vehicle active suspension systems, more sophisticated reward formulations have been developed to handle higher-dimensional dynamics and partially observable conditions, highlighting the importance of problem-specific reward design in practical applications [38]. Future work will explore integrating MPC to guide long-horizon planning and improve sample efficiency.

In addition, while this work demonstrates cross-generation performance improvement within the proposed framework, a rigorous comparison against established uncertainty-aware CCD methods (e.g., robust and stochastic MPC approaches) and formal evaluation of uncertainty calibration metrics (such as quantile coverage and reliability) remain important directions for future investigation. Another future work will investigate more case studies of active suspension systems across diverse scenarios, aiming to optimize designs for different vehicle types, driving styles, and regional conditions. Given that optimal system configurations vary based on operational requirements, our proposed framework's ability to autonomously learn from the environment and cooptimize both physical components and control policies is particularly valuable. One key direction worth investigating is how we achieve optimal designs for varying driving speed profiles and road conditions. Furthermore, while this work focuses on active suspension, the methodology is generalizable to other engineering applications, such as wind turbines and autonomous vehicles, where dynamic interactions between control and physical system design play a crucial role.

## Acknowledgment

We are grateful for the grant support from the National Science Foundation's (NSF) Engineering Research Center for Hybrid Autonomous Manufacturing: Moving from Evolution to Revolution (ERC-HAMMER) and NSF's FMRG: Manufacturing ADvanced Electronics through Printing Using Bio-based and Locally Identifiable Compounds (MADE-PUBLIC), under

Award Numbers EEC-2133630 and CMMI-2037026, respectively. Yi-Ping Chen also acknowledges the Taiwan-Northwestern Doctoral Scholarship to support his doctoral study. We also thank Dr. Anton van Beek at University College Dublin for his feedback and insightful discussions.

## Conflict of Interest

There are no conflicts of interest.

## Data Availability Statement

The datasets generated and supporting the findings of this article are obtainable from the corresponding author upon reasonable request.

## References

- [1] Garcia-Sanz, M., 2019, "Control Co-Design: An Engineering Game Changer," *Adv. Control Appl.: Eng. Ind. Syst.*, **1**(1), p. e18.
- [2] Allison, J. T., and Herber, D. R., 2014, "Special Section on Multidisciplinary Design Optimization: Multidisciplinary Design Optimization of Dynamic Engineering Systems," *AIAA J.*, **52**(4), pp. 691–710.
- [3] Tsai, Y.-K., Malak, J., and Richard, J., 2022, "A Constraint-Handling Technique for Parametric Optimization and Control Co-Design," International Design Engineering Technical Conferences and Computers and Information in Engineering Conference, St. Louis, MO, Aug. 14–17.
- [4] Allison, J. T., Guo, T., and Han, Z., 2014, "Co-Design of an Active Suspension Using Simultaneous Dynamic Optimization," *ASME J. Mech. Des.*, **136**(8), p. 081003.
- [5] Bayat, S., and Allison, J. T., 2026, "Control Co-Design With Varying Available Information Applied to Vehicle Suspensions," *ASME J. Dyn. Syst., Meas., Control.*, **148**(1), p. 011013.
- [6] Tsai, Y.-K., and Malak, R. J., 2026, "Parametric Multi-Objective Optimization for Simultaneous and Nested Control Co-Design Formulations With Tube-Based Model Predictive Controllers," *ASME J. Mech. Des.*, **148**(11), p. 111701.
- [7] Tsai, Y.-K., 2023, "Control Co-Design Using Parametric Optimization," Ph.D. Dissertation, Texas A&M University, College Station, TX.
- [8] Bravo-Palacios, G., Grandesso, G., Prete, A. D., and Wensing, P. M., 2022, "Robust Co-Design: Coupling Morphology and Feedback Design Through Stochastic Programming," *ASME J. Dyn. Syst. Meas. Control.*, **144**(2), p. 021007.
- [9] Nash, A. L., Pangborn, H. C., and Jain, N., 2021, "Robust Control Co-Design With Receding-Horizon MPC," 2021 American Control Conference (ACC), New Orleans, LA, May 25–28, IEEE, pp. 373–379.
- [10] Tsai, Y.-K., and Malak, R. J., 2023, "Robust Control Co-Design Using Tube-Based Model Predictive Control," 2023 American Control Conference (ACC), San Diego, CA, May 31–June 2, IEEE, pp. 769–775.
- [11] Tsai, Y.-K., and Malak, R., 2025, "Control Co-Design With Performance-Robustness Trade-Off Using Tube-Based Stochastic Model Predictive Control," *ASME Lett. Dyn. Syst. Control.*, **5**(3), p. 030903.
- [12] Cui, T., Allison, J. T., and Wang, P., 2021, "Reliability-Based Control Co-Design of Horizontal Axis Wind Turbines," *Struct. Multidiscipl. Optim.*, **64**(6), pp. 3653–3679.
- [13] Rudnick-Cohen, E. S., Hodson, J. D., Reich, G. W., Pankonien, A. M., and Beran, P. S., 2022, "Robust Optimal Design and Control of a Maneuvering Morphing Airfoil," *J. Aircr.*, **59**(4), pp. 861–874.
- [14] Azad, S., and Herber, D. R., 2023, "An Overview of Uncertain Control Co-Design Formulations," *ASME J. Mech. Des.*, **145**(9), p. 091709.
- [15] Fine, J. B., Holbrook, I., and Vermillion, C., 2024, "Co-Design for Real-Time Adaptability: Methodology and Wind Energy Case Study," *IFAC-PapersOnLine*, **58**(28), pp. 816–821.
- [16] Sundararajan, A., and Herber, D. R., 2023, "Using High-Fidelity Time-Domain Simulation Data to Construct Multi-Fidelity State Derivative Function Surrogate Models for Use in Control and Optimization," ASME International Mechanical Engineering Congress and Exposition, Vol. 87639, American Society of Mechanical Engineers, p. V006T07A089.
- [17] Chen, Y.-P., and Chan, K.-Y., 2021, "Unknown Parameter Excitation and Estimation for Complex Systems With Dynamic Performances," *ASME J. Mech. Des.*, **143**(9), p. 091704.
- [18] Ren, L., and Xi, Z., 2022, "Bias-Learning-Based Model Predictive Controller Design for Reliable Path Tracking of Autonomous Vehicles Under Model and Environmental Uncertainty," *ASME J. Mech. Des.*, **144**(9), p. 091706.
- [19] Chen, T., He, Z., and Ciocarlie, M., 2021, "Co-Designing Hardware and Control for Robot Hands," *Sci. Rob.*, **6**(54), p. eabg2133.
- [20] Sadat, E., Lotfalian Saremi, M., and Bayrak, A. E., 2023, "A Two-Timescale Reinforcement Learning Approach for Control Co-Design Problems," International Design Engineering Technical Conferences and Computers and Information in Engineering Conference, Vol. 87301, American Society of Mechanical Engineers, p. V03AT03A003.

- [21] Chen, C., Xiang, P., Zhang, J., Xiong, R., Wang, Y., and Lu, H., 2023, "Deep Reinforcement Learning Based Co-Optimization of Morphology and Gait for Small-Scale Legged Robot," *IEEE/ASME Trans. Mechatron.*, **29**(4), pp. 2697–2708.
- [22] Chen, X., Huang, D., Li, M., Cai, Y., Wen, Z., Cai, Z., and Yang, W., 2023, "Evolving Physical Instinct for Morphology and Control Co-Adaption," 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), IEEE, pp. 6616–6623.
- [23] Chen, T., He, Z., and Ciocarlie, M., 2020, "Hardware as Policy: Mechanical and Computational Co-Optimization Using Deep Reinforcement Learning," arXiv preprint [arXiv:2008.04460](https://arxiv.org/abs/2008.04460).
- [24] Sun, J., Yao, M., Xiao, X., Xie, Z., and Zheng, B., 2023, "Co-Optimization of Morphology and Behavior of Modular Robots via Hierarchical Deep Reinforcement Learning," Robotics: Science and Systems, Daegu, Republic of Korea, July 10–14.
- [25] Schaff, C., Yunis, D., Chakrabarti, A., and Walter, M. R., 2019, "Jointly Learning to Construct and Control Agents Using Deep Reinforcement Learning," 2019 International Conference on Robotics and Automation (ICRA), IEEE, pp. 9798–9805.
- [26] Luck, K. S., Amor, H. B., and Calandra, R., 2020, "Data-Efficient Co-Adaptation of Morphology and Behaviour With Deep Reinforcement Learning," Conference on Robot Learning, PMLR, pp. 854–869.
- [27] Karkaria, V., Tsai, Y.-K., Chen, Y.-P., and Chen, W., 2025, "An Optimization-Centric Review for Integrating Artificial Intelligence and Digital Twin Technologies in Manufacturing," *Eng. Optim.*, **57**(1), pp. 161–207.
- [28] Willcox, K., Bingham, D., Chung, C., Chung, J., Cruz-Neira, C., Grant, C., Kinter, J., Leung, R., et al., 2023, *Foundational Research Gaps and Future Directions for Digital Twins*, Technical Report. National Academies of Sciences, Engineering, and Medicine.
- [29] Thelen, A., Zhang, X., Fink, O., Lu, Y., Ghosh, S., Youn, B. D., Todd, M. D., Mahadevan, S., Hu, C., and Hu, Z., 2022, "A Comprehensive Review of Digital Twin—Part 1: Modeling and Twinning Enabling Technologies," *Struct. Multidiscipl. Optim.*, **65**(12), p. 354.
- [30] Thelen, A., Zhang, X., Fink, O., Lu, Y., Ghosh, S., Youn, B. D., Todd, M. D., Mahadevan, S., Hu, C., and Hu, Z., 2023, "A Comprehensive Review of Digital Twin—Part 2: Roles of Uncertainty Quantification and Optimization, a Battery Digital Twin, and Perspectives," *Struct. Multidiscipl. Optim.*, **66**(1), p. 1.
- [31] Zemskov, A. D., Fu, Y., Li, R., Wang, X., Karkaria, V., Tsai, Y.-K., Chen, W., et al., 2024, "Security and Privacy of Digital Twins for Advanced Manufacturing: A Survey," arXiv preprint [arXiv:2412.13939](https://arxiv.org/abs/2412.13939)
- [32] Karkaria, V., Goeckner, A., Zha, R., Chen, J., Zhang, J., Zhu, Q., Cao, J., Gao, R. X., and Chen, W., 2024, "Towards a Digital Twin Framework in Additive Manufacturing: Machine Learning and Bayesian Optimization for Time Series Process Optimization," *J. Manuf. Syst.*, **75**, pp. 322–332.
- [33] van Beek, A., Neville Karkaria, V., and Chen, W., 2023, "Digital Twins for the Designs of Systems: A Perspective," *Struct. Multidiscipl. Optim.*, **66**(3), p. 49.
- [34] Chen, Y.-P., Karkaria, V., Tsai, Y.-K., Rolark, F., Quispe, D., Gao, R. X., Cao, J., and Chen, W., 2025, "Real-Time Decision-Making for Digital Twin in Additive Manufacturing With Model Predictive Control Using Time-Series Deep Neural Networks," *J. Manuf. Syst.*, **80**, pp. 412–424.
- [35] Bhatti, G., Mohan, H., and Singh, R. R., 2021, "Towards the Future of Smart Electric Vehicles: Digital Twin Technology," *Renew. Sustain. Energy Rev.*, **141**, p. 110801.
- [36] Schena, L., Marques, P. A., Poletti, R., Ahizi, S., Van den Berghe, J., and Mendez, M. A., 2024, "Reinforcement Twinning: From Digital Twins to Model-Based Reinforcement Learning," *J. Comput. Sci.*, **82**, p. 102421.
- [37] Xia, K., Sacco, C., Kirkpatrick, M., Saidu, C., Nguyen, L., Kircaliali, A., and Harik, R., 2021, "A Digital Twin to Train Deep Reinforcement Learning Agent for Smart Manufacturing Plants: Environment, Interfaces and Intelligence," *J. Manuf. Syst.*, **58**, pp. 210–230.
- [38] Tsai, Y.-K., Chen, Y.-P., Karkaria, V., and Chen, W., 2026, "Reinforcement Learning-Based Control Co-Design of Digital Twin-Enabled Full-Vehicle Active Suspension Systems," *Struct. Multidiscipl. Optim.*, **69**.
- [39] Herber, D. R., and Allison, J. T., 2019, "Nested and Simultaneous Solution Strategies for General Combined Plant and Control Design Problems," *ASME J. Mech. Des.*, **141**(1), p. 011402.
- [40] Buşoniu, L., De Bruin, T., Tolić, D., Kober, J., and Palunko, I., 2018, "Reinforcement Learning for Control: Performance, Stability, and Deep Approximators," *Ann. Rev. Control.*, **46**, pp. 8–28.
- [41] Sutton, R. S., and Barto, A. G., 1998, *Reinforcement Learning: An Introduction*, Vol. 1, MIT Press, Cambridge.
- [42] Shakya, A. K., Pillai, G., and Chakrabarty, S., 2023, "Reinforcement Learning Algorithms: A Brief Survey," *Expert Syst. Appl.*, **231**, p. 120495.
- [43] Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P., 2015, "Trust Region Policy Optimization," International Conference on Machine Learning, PMLR, pp. 1889–1897.
- [44] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O., 2017, "Proximal Policy Optimization Algorithms," arXiv preprint [arXiv:1707.06347](https://arxiv.org/abs/1707.06347)
- [45] Arulkumaran, K., Deisenroth, M. P., Brundage, M., and Bharath, A. A., 2017, "Deep Reinforcement Learning: A Brief Survey," *IEEE Signal Process. Mag.*, **34**(6), pp. 26–38.
- [46] Ladosz, P., Weng, L., Kim, M., and Oh, H., 2022, "Exploration in Deep Reinforcement Learning: A Survey," *Inf. Fusion*, **85**, pp. 1–22.
- [47] Wang, X., Wang, S., Liang, X., Zhao, D., Huang, J., Xu, X., Dai, B., and Miao, Q., 2022, "Deep Reinforcement Learning: A Survey," *IEEE Trans. Neural Netw. Learn. Syst.*, **35**(4), pp. 5064–5078.
- [48] Ma, P., Du, T., Zhang, J. Z., Wu, K., Spielberg, A., Katschmann, R. K., and Matusik, W., 2021, "DiffAqua: A Differentiable Computational Design Pipeline for Soft Underwater Swimmers With Shape Interpolation," *ACM Trans. Graph.*, **40**(4), pp. 1–14.
- [49] Yuhn, C., Sato, Y., Kobayashi, H., Kawamoto, A., and Nomura, T., 2023, "4D Topology Optimization: Integrated Optimization of the Structure and Self-Actuation of Soft Bodies for Dynamic Motions," *Comput. Methods Appl. Mech. Eng.*, **414**, p. 116187.
- [50] Wang, Y., Wu, S., Zhang, T., Chang, Y., Fu, H., Fu, Q., and Wang, X., 2023, "Preco: Enhancing Generalization in Co-Design of Modular Soft Robots via Brain-Body Pre-Training," Conference on Robot Learning, PMLR, pp. 478–498.
- [51] Wang, T., Zhou, Y., Fidler, S., and Ba, J., 2019, "Neural Graph Evolution: Towards Efficient Automatic Robot Design," International Conference on Learning Representations, New Orleans, LA, May 6–9.
- [52] Gupta, A., Savarese, S., Ganguli, S., and Fei-Fei, L., 2021, "Embodied Intelligence via Learning and Evolution," *Nat. Commun.*, **12**(1), p. 5721.
- [53] Chen, C., Xiang, P., Lu, H., Wang, Y., and Xiong, R., 2023, "C 2: Co-Design of Robots via Concurrent-Network Coupling Online and Offline Reinforcement Learning," 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), IEEE, pp. 7487–7494.
- [54] Mayne, D. Q., Seron, M. M., and Raković, S. V., 2005, "Robust Model Predictive Control of Constrained Linear Systems With Bounded Disturbances," *Automatica*, **41**(2), pp. 219–224.
- [55] Chen, S., Saulnier, K., Atanasov, N., Lee, D. D., Kumar, V., Pappas, G. J., and Morari, M., 2018, "Approximating Explicit Model Predictive Control Using Constrained Neural Networks," 2018 Annual American Control Conference (ACC)IEEE, pp. 1520–1527.
- [56] Tsai, Y.-K., and Malak, R. J., 2023, "Control Co-Design With Approximate Explicit Model Predictive Controllers," International Design Engineering Technical Conferences and Computers and Information in Engineering Conference, Boston, MA, Aug. 20–23.
- [57] Herceg, M., Kvasnica, M., Jones, C., and Morari, M., 2013, "Multi-Parametric Toolbox 3.0," Proceedings of the European Control Conference, Zürich, Switzerland, pp. 502–510. <http://control.ee.ethz.ch/~mpt>
- [58] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., et al., 2015, "Human-Level Control Through Deep Reinforcement Learning," *Nature*, **518**(7540), pp. 529–533.
- [59] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., and Killeen, T., 2019, "Pytorch: An Imperative Style, High-Performance Deep Learning Library," NeurIPS 2019, Vancouver, BC, Canada, Dec. 8–14.
- [60] Arendt, P. D., Apley, D. W., and Chen, W., 2012, "Quantification of Model Uncertainty: Calibration, Model Discrepancy, and Identifiability," *ASME J. Mech. Des.*, **134**(10), p. 100908.
- [61] Chen, Y.-P., Tsai, Y.-K., Karkaria, V., and Chen, W., 2026, "Uncertainty-Aware Digital Twins: Robust Model Predictive Control Using Time-Series Deep Quantile Learning," *ASME J. Mech. Des.*, **148**(2), p. 021702.
- [62] Chen, Y.-P., Suarez, D., Tsai, Y.-K., Karkaria, V., Hu, G., Chen, Z., Guo, P., Cao, J., and Chen, W., 2026, "Adaptive Digital Twin of Sheet Metal Forming Via Proper Orthogonal Decomposition-Based Koopman Operator With Model Predictive Control," *J. Manuf. Syst.*, **86**, pp. 762–773.