# Contraction Actor-Critic:
# Contraction Metric-Guided Reinforcement Learning for Robust Path Tracking

**Minjae Cho**
The Grainger College of Engineering
University of Illinois Urbana-Champaign, US
`minjae5@illinois.edu`

**Hiroyasu Tsukamoto**
The Grainger College of Engineering
University of Illinois Urbana-Champaign, US
`hiroyasu@illinois.edu`

**Huy T. Tran**
The Grainger College of Engineering
University of Illinois Urbana-Champaign, US
`huytran1@illinois.edu`

**Abstract:** Control contraction metrics (CCMs) provide a framework to co-synthesize a controller and a corresponding contraction metric—a positive-definite Riemannian metric under which a closed-loop system is guaranteed to be incrementally exponentially stable. However, the synthesized controller only ensures that all the trajectories of the system converge to one single trajectory and, as such, does not impose any notion of optimality across an entire trajectory. Furthermore, constructing CCMs requires a known dynamics model and non-trivial effort in solving an infinite-dimensional convex feasibility problem, which limits its scalability to complex systems featuring high dimensionality with uncertainty. To address these issues, we propose to integrate CCMs into reinforcement learning (RL), where CCMs provide dynamics-informed feedback for learning control policies that minimize cumulative tracking error under unknown dynamics. We show that our algorithm, called contraction actor-critic (CAC), formally enhances the capability of CCMs to provide a set of contracting policies with the long-term optimality of RL in a fully automated setting. Given a pre-trained dynamics model, CAC simultaneously learns a contraction metric generator (CMG)—which generates a contraction metric—and uses an actor-critic algorithm to learn an optimal tracking policy guided by that metric. We demonstrate the effectiveness of our algorithm relative to established baselines through extensive empirical studies, including simulated and real-world robot experiments, and provide a theoretical rationale for incorporating contraction theory into RL.

**Keywords:** Contraction Theory, Control Contraction Metric, Reinforcement Learning, Path-Tracking Control

## 1 Introduction

A fundamental challenge in designing path-tracking controllers, particularly those in safety-critical systems, is to ensure safe and reliable operation while simultaneously achieving task objectives. Control contraction metrics (CCMs) [1], grounded in contraction theory [2], provide a formal approach for certifying the incremental exponential stability of a path-tracking controller with a desired reference trajectory. This certificate is contingent upon the existence of a controller and a corresponding contraction metric—a positive-definite Riemannian metric—that satisfies the contraction and CCM conditions. However, this certificate does not guarantee that the resulting controller optimally minimizes the path-tracking error across the entire trajectory. Moreover, existing

work [3, 4, 5, 6] requires a known dynamics model and non-trivial effort in solving an infinite-dimensional convex feasibility problem, limiting its scalability to complex systems featuring high dimensionality with uncertainty.

[7] proposes a learning framework that scales existing CCM methods to complex systems, including a 10-dimensional nonlinear system of a quadrotor. Their framework jointly trains two neural networks—one for the controller and another for a contraction metric generator (CMG) that generates a contraction metric—to minimize violations of the contraction and CCM conditions with the contraction rate treated as a fixed hyperparameter. However, this setup lacks any notion of optimality, even myopically, with respect to cumulative tracking error and relies on known dynamics to evaluate the satisfaction of contraction and CCM conditions. While function approximation can be used to learn the required dynamics, [8] shows that the method proposed in [7] is ineffective in learning a contracting policy when the dynamics are approximated.

To this end, we propose to integrate CCMs into reinforcement learning (RL) to optimize control policies that minimize cumulative tracking error in complex nonlinear systems with unknown dynamics. Our algorithm, contraction actor-critic (CAC), integrates the optimality and robustness benefits of model-free RL with CCMs to provide a contraction certificate in a fully automated manner. Given a pre-trained dynamics model, CAC simultaneously uses that model to learn a CMG and uses an actor-critic algorithm to learn a policy guided by the CMG. We demonstrate the performance improvement of our algorithm compared to established baselines through extensive empirical studies, including simulated and real-world robot experiments, and provide a theoretical analysis of our approach.

## 2 Preliminaries

For a closed-loop system, we denote a policy as $\pi : \mathcal{X} \times \mathcal{U} \to [0, 1]$, where $\mathcal{X} \subseteq \mathbb{R}^n$ is an $n$-dimensional compact state set and $\mathcal{U} \subseteq \mathbb{R}^m$ is an $m$-dimensional compact control set. We then consider the following control-affine system,

$$\dot{x}(t) = f(x(t)) + B(x(t))u(t), \tag{1}$$

where $f : \mathbb{R}^n \to \mathbb{R}^n$ denotes the smooth nonlinear drift dynamics and $B : \mathbb{R}^n \to \mathbb{R}^{n \times m}$ is the smooth actuation matrix function. For notational simplicity, we omit the explicit time dependence to write $x(t)$ and $u(t)$ simply as $x$ and $u$, unless the time argument is required for clarity.

### 2.1 Contraction Theory and Control Contraction Metrics

Contraction theory [2, 9, 4] provides a framework for analyzing the convergence between trajectories that satisfy Equation (1) by leveraging the differential form of Lyapunov functions. Specifically, contraction theory reformulates conventional Lyapunov stability conditions by employing a quadratic Lyapunov function defined over differential dynamics,

$$V(x(t), \delta x(t)) = \delta x^\top M(x(t)) \delta x, \tag{2}$$

where $\delta x(t)$ represents the infinitesimal displacement between two trajectories at time $t$ and $M(x(t)) \succ 0$ is a symmetric, positive definite contraction metric. When the time derivative of $V$ satisfies $\dot{V} \leq -2\lambda V$ for some $\lambda > 0$, the differential dynamics, $\delta x(t)$, are guaranteed to contract exponentially at rate $\lambda$.

**Control contraction metrics.** [1] extends contraction theory by integrating the effect of control inputs. Formally, the system is said to be contracting under a policy $\pi$ if there exists a metric $M(x) \succ 0$ and a control policy $u \sim \pi(x)$ such that,

$$\dot{M} + \text{sym}\left(M(A + BK)\right) + 2\lambda M \preceq 0, \tag{3}$$

where $\text{sym}(\cdot)$ is the symmetric part of the input argument (i.e., $\text{sym}(A) = A + A^\top$), $A(x, u) = \frac{\partial f}{\partial x} + \sum_{i=1}^m u_i \frac{\partial b_i}{\partial x}$ is the Jacobian of the drift dynamics, $K(x, t) = \frac{\partial u}{\partial x}$ represents the feedback gain associated with the policy, and $\lambda > 0$ specifies the contraction rate under the metric $M$. Satisfaction

of this condition ensures that the policy is guaranteed to exhibit incremental exponential convergence to the reference trajectory, serving as a strong certificate of stability in a time-varying system.

A sufficient condition for the existence of a feedback control policy satisfying Equation (3) is given by the following constraints,

$$B_\perp^\top \left( -\partial_f W(x) + \text{sym}\left( \frac{\partial W(x)}{\partial x} f(x) \right) + 2\lambda W(x) \right) B_\perp \prec 0, \tag{4}$$

$$B_\perp^\top \left( \partial_{b_j} W(x) - \text{sym}\left( \frac{\partial b_j(x)}{\partial x} W(x) \right) \right) B_\perp = 0, \quad j = 1, \dots, m. \tag{5}$$

Here, $W$ denotes the dual metric of $M$ with $W = M^{-1}$, $\partial_f W(x)$ represents the Lie derivative of the matrix-valued function $W(x)$ along the vector $f \in \mathbb{R}^n$, defined as $\partial_f W := \sum_{i=1}^n f_i \frac{\partial W}{\partial x_i}$, and $B_\perp$ denotes the null space of the actuation matrix $B^\top$. The constraints above ensure that the dynamics that cannot be directly manipulated by the policy must naturally contract. The Riemannian metric defined with the positive definite matrix $M$ of Equation (4) and Equation (5) is called the control contraction metric.

## 2.2 Reinforcement Learning

RL addresses decision-making under uncertainty, where an agent learns to make decisions by interacting with an environment, typically formalized as a Markov decision process (MDP) [10].

**Markov decision process.** Building on the previously defined state and control sets, an MDP is defined by a tuple $\langle \mathcal{X}, \mathcal{U}, T, R, \gamma \rangle$, where $T : \mathcal{X} \times \mathcal{U} \times \mathcal{X} \to [0, 1]$ defines the probability of transitioning to state $x'$ given the current state $x$ and control input $u$, and $R : \mathcal{X} \times \mathcal{U} \times \mathcal{X} \to \mathbb{R}$ is the reward function. The discount factor $\gamma \in (0, 1]$ weighs future rewards relative to immediate ones, encouraging the agent to consider long-term consequences. The goal of the agent is to learn a control policy $\pi$ that maximizes the expected cumulative rewards. The cumulative rewards (return) at time $t$ is given by $G_t = \sum_{k=0}^\infty \gamma^k r_{t+k+1}$, where $r_t$ is the reward received at time step $t$.

**Actor-critic algorithms.** Actor-critic algorithms [11, 12, 13, 14] integrate policy-based and value-based learning in RL. The *actor*, represented by a policy $\pi_\theta$ parameterized by $\theta$, gives control inputs for each state. Concurrently, the *critic* estimates the value function $\mathcal{V}_\phi^{\pi_\theta}(x) = \mathbb{E}_{\pi_\theta}[G_t \mid x_t = x]$ for policy $\pi_\theta$, parameterized by $\phi$ [1]. The critic is used to estimate the policy gradient used to refine the policy. Specifically, the policy parameters are updated via policy gradient methods [15] as follows,

$$\theta^{i+1} = \theta^i + \nabla_{\theta^i} J(\theta^i), \tag{6}$$

where $\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(u_t \mid x_t) \cdot G_t \right]$ is an estimated policy gradient and $i$ denotes the iteration number. The critic is trained to minimize the Bellman error,

$$\phi = \underset{\phi}{\arg\min} \, \mathbb{E}_{\pi_\theta} \left[ \left( \mathcal{V}_\phi(x_t) - (r_t + \gamma \mathcal{V}_\phi(x_{t+1})) \right)^2 \right]. \tag{7}$$

## 3 Our Algorithm: Contraction Actor-Critic (CAC)

We incorporate contraction certificates into an actor-critic framework to synthesize an effective and robust path-tracking policy under unknown dynamics. Our approach is illustrated in Figure 1. We first pre-train a model to approximate the system dynamics. We then simultaneously learn a CMG and a policy guided by the CMG. Intuitively, the CMG finds a contraction metric, which we then use to define a reward function for the policy that incentivizes it to achieve a high contraction rate. Our resulting policies minimize cumulative tracking error while also inheriting a contraction certificate. While we use an actor-critic algorithm, in principle, our approach could be used with any RL algorithm. Note that the learned dynamics are only used to evaluate the contraction and CCM conditions in the CMG. Our algorithm is outlined in Algorithm 1 and described in detail below.

---

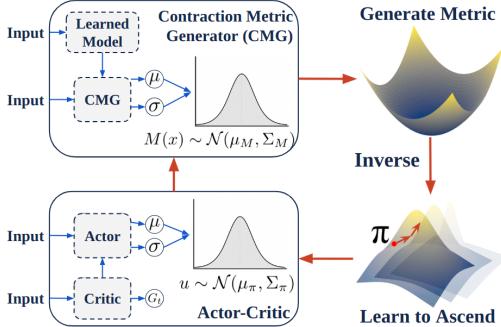[1]We use $\mathcal{V}$ for value function to differentiate it from a Lyapunov function $V$.

Figure 1: We jointly train a CMG and a policy using RL, where the reward is defined by the CMG in a manner that ensures contraction.

**Algorithm 1 CAC**: Contraction Actor-Critic

**Require:** Data: $\mathcal{D} = \{\dot{x}_i, x_i, u_i\}_{i=0}^N$
**Require:** policy iterations $n$
1: Initialize $\xi, \zeta, \chi, \theta$, and $\phi$
2: /* Dynamics pre-training */
3: **for** each batch $d$ from $\mathcal{D}$ **do**
4:      Update $\xi, \zeta$ with Eqn. 8
5: **end for**
6: /* CAC training */
7: **for** each batch $d$ from $\pi, T$ **do**
8:      Update $\chi$ with Eqn. 9 *(every n-th step)*
9:      Update $\theta$ and $\phi$ with Eqn. 6 and 7
10: **end for**

### 3.1 Learning a Dynamics Model

To evaluate the contraction and CCM conditions given by Equations (3) to (5), we need a model of the drift dynamics $f$, the actuation matrix function $B$, and its null space $B_\perp$. Following [8], we employ neural network function approximators to model $f$ and $B$ and use numerical techniques to compute $B_\perp$. Specifically, given a dataset $\mathcal{D} = \{(\dot{x}_i, x_i, u_i)\}_{i=0}^N$, we simultaneously learn two separate neural networks, $\hat{f}_\xi$ and $\hat{B}_\zeta$, parameterized by $\xi$ and $\zeta$ respectively, as follows,

$$\xi, \zeta = \arg\min_{\xi, \zeta} \mathbb{E}_{(\dot{x}, x, u) \sim \mathcal{D}} \left[ \left\| \dot{x} - \left( \hat{f}_\xi(x) + \hat{B}_\zeta(x)u \right) \right\|_2^2 \right]. \tag{8}$$

We directly compute $B_\perp$ using singular value decomposition (SVD), where we select the left-singular vectors beyond the rank of $B$. More specifically, since $\text{SVD}(B) = U\Sigma V^\top$ and $B^\top = V\Sigma^\top U^\top$, the null space of $B^\top$ is spanned by the columns of $U$ associated with the indices beyond the rank of $B$. For example, if $B \in \mathbb{R}^{3 \times 2}$ has rank 2, then $U \in \mathbb{R}^{3 \times 3}$ and the third column of $U$ spans the null space of $B^\top$.

### 3.2 Jointly Learning a Contraction Metric Generator and Policy

To enable contraction-guided RL, we simultaneously learn a CMG and a policy maximizing rewards defined by that CMG. However, this process can be unstable because the CMG is conditioned on the current policy. We address this issue by implementing a freeze-and-learn strategy, where for each update of the CMG, we freeze its parameters and then perform $n$ policy updates to allow the RL policy to adapt to the newly defined reward landscape.

**Learning a contraction metric generator.** Given a learned dynamics model, we train the CMG to output a Gaussian distribution over contraction metrics. Specifically, we use a probabilistic model to enable the use of entropy regularization, which has been shown to result in faster and more stable convergence for policy optimization [16]. This approach also prevents premature convergence of the CMG in our bi-level optimization setting. Specifically, we learn a CMG, $M \sim \mathcal{M}_\chi(x)$ parameterized by $\chi$, based on the loss used in [7] with an additional reward-conditioned entropy regularizer. Specifically, let $C_M$, $C_{W_1}$, and $C_{W_2}^j$ denote the left-hand side of Equation (3), the inequality constraint in Equation (4), and the $j$th equality constraint in Equation (5) for $j = 1, \ldots, m$, respectively. Then the CMG is learned as follows,

$$\chi = \underset{\chi}{\arg\min} \, \mathbb{E}_{\pi_\theta} \left[ L_{PD}(-C_M) + L_{PD}(-C_{W_1}) + \sum_{j=1}^m \|C_{W_2}^j\|_F - \alpha(r_t)\mathcal{H}(\mathcal{M}_\chi(x_t)) \right], \tag{9}$$

where $L_{PD}(A) = z^\top A z$ if $A \prec 0$ else $0$ is a loss that penalizes non-positiveness of the input matrix with $z \sim U(-1, 1)$, $\|\cdot\|_F$ is a Frobenius matrix norm, and $\mathcal{H}(\mathcal{M}_\chi(x_t))$ is the entropy of the

4

probability distribution $\mathcal{M}_\chi$ for state $x_t$. The strength of the entropy regularizer is regulated by the reward, $r_t$, achieved by the policy such that $\alpha(r_t) = \beta_\mathcal{M} e^{-r_t}$ where $\beta_\mathcal{M} > 0$ is a hyperparameter. Intuitively, the reward-conditioned entropy regularizer encourages the CMG to explore the CCM space when the agent receives low returns and to exploit when high returns are achieved.

**Learning a contraction-guided policy.** Suppose we have a contraction metric, $M \sim \mathcal{M}_\chi(x)$, that satisfies the contraction and CCM conditions in Equations (3) to (5). This guarantees that there exists a set of control inputs $\tilde{\mathcal{U}} \subseteq \mathcal{U}$ that ensures system contraction [1, 17, 4]. While a control $u \in \tilde{\mathcal{U}}$ satisfies the contraction condition under the given metric, it is not guaranteed to minimize the cumulative tracking error. We therefore use RL to optimize a policy that does so, by defining a reward function based on the differential Lyapunov function given by Equation (2). Intuitively, this reward encourages the policy to select controls that induce the optimal contraction within the admissible set $\tilde{\mathcal{U}}$.

We specifically use the following reward function,

$$R(x) = \frac{1}{1 + \delta x^\top M \delta x} + \beta_\pi \mathcal{H}(\pi_\theta(x)), \tag{10}$$

where the state $x$ includes the current time step $t$, $\delta x$ is the infinitesimal displacement between $x$ and the reference trajectory at time $t$, $M \sim \mathcal{M}_\chi(x)$ is a contraction metric, and $\mathcal{H}(\pi_\theta(x))$ is an entropy regularizer with scaler $\beta_\pi > 0$. In practice, we calculate the infinitesimal displacement as the Euclidean displacement between the two points. We include the one in the denominator of $1/(1 + \delta x^\top M \delta x)$ to ensure that term is bounded within $[0, 1]$, motivated by prior work that has shown that bounded rewards can improve numerical stability [18, 19, 20].

## 4 Theoretical Results: Bridging RL and Contraction Theory

We now give theoretical justification for our idea to integrate a contraction metric into an RL reward function. Specifically, we prove asymptotic convergence of tracking errors by an RL policy trained to maximize Equation (10) without entropy regularization (i.e., $\beta_\pi = 0$).

Assume that there exists a contracting policy $\pi_c$ which satisfies the contraction and CCM conditions.

**Assumption 1.** *There exists at least one contracting policy $\pi_c$ that, along with contraction metric $M$, satisfies Equations (3) to (5) with a rate of $\lambda > 0$ such that the policy $\pi_c$ satisfies,*

$$\|\delta x(t)\|_M^2 \le \|\delta x(t_0)\|_M^2 e^{-2\lambda(t-t_0)}, \tag{11}$$

where $t$ is the current time, $t_0$ is the initial time, and $\|\delta x(t)\|_M^2 = \int_0^1 \frac{\partial q}{\partial \mu}^\top M(q(\mu, t)) \frac{\partial q}{\partial \mu} d\mu$ denotes the integration over the geodesics connecting $x(t)$ and $x_d(t)$, with $q(\mu = 0, t) = x(t)$ and $q(\mu = 1, t) = x_d(t)$.

Now, define the performance of a contracting policy over an infinite horizon $[t_0, \infty]$ as follows,

$$\mathcal{J}_T^{\pi_c}(x_0, t_0) := \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \|\delta x(t_k)\|_M^2 \right], \tag{12}$$

where $x_0 = x(t_0)$ is the initial state and $\Delta t$ is the discrete time interval such that the current time at $k^{th}$ sample is given by $t_k = k\Delta t + t_0$. Based on Equation (11), we derive the following lemma.

**Lemma 1** (Bound of Contraction Inequality). *If* (11) *holds, the performance measure in Equation* (12) *satisfies the following upper bound,*

$$\mathcal{J}_T^{\pi_c}(x_0, t_0) = \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \|\delta x(t_k)\|_M^2 \right] \le \frac{\|\delta x(t_0)\|_M^2}{1 - e^{2\lambda \Delta t}}.$$

*Proof.* Can be found in Appendix A.1. □

Consider a reward function defined by Equation (10) but without an entropy term,

$$\tilde{R}(x) = \frac{1}{1 + \|\delta x\|_M^2} \in [0, 1]. \tag{13}$$

We define a performance measure for an arbitrary policy $\pi$ and its RL-trained optimal counterpart $\pi^*$ as follows,

$$\tilde{\mathcal{J}}_T^\pi(x_0, t_0) := \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k \tilde{R}_k \right], \quad \tilde{\mathcal{J}}_T^{\pi^*}(x_0, t_0) := \max_\pi \tilde{\mathcal{J}}_T^\pi(x_0, t_0), \tag{14}$$

where $\tilde{R}_k := \tilde{R}(x(t_k))$ for notational clarity.

We then introduce a cost function to establish equivalence between minimizing that cost and maximizing the performance measure defined in Equation (14).

**Lemma 2** (Equivalence between Cost Minimization and Reward Maximization)**.** *Define a cost function in terms of $\tilde{R}(x)$ as,*

$$C(x) = 1 - \tilde{R}(x) \in [0, 1].$$

*Then, the following equivalence holds,*

$$\pi^* = \operatorname*{argmax}_\pi \tilde{\mathcal{J}}_T^\pi(x_0, t_0) = \operatorname*{argmax}_\pi \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k \tilde{R}_k \right] = \operatorname*{argmin}_\pi \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k C_k \right].$$

where $C_k := C(x(t_k))$ for notational clarity.

*Proof.* Can be found in Appendix A.2. $\qquad\square$

Finally, using Lemmas 1 and 2, we now prove asymptotic convergence of an RL policy trained to maximize the objective given in Equation (14), given that Assumption 1 is satisfied.

**Theorem 1** (Asymptotic Convergence of an RL-trained Policy via Existence of a Contracting Policy)**.** *If there exists at least one contracting policy $\pi_c$ with a contraction rate $\alpha > 0$, then the optimal policy $\pi^*$ obtained via Equation (14) must exhibit asymptotic convergence, as given by,*

$$\lim_{t \to \infty} \|\delta x(t)\|_M^2 = 0, \ \forall x_0, t_0.$$

*Proof.* Can be found in Appendix A.3. $\qquad\square$

## 5 Experimental Results

We conduct experiments in both simulated and real-world settings to validate the effectiveness and efficiency of our method under model approximation errors and sim-to-real dynamics deviations. All of our results assume unknown dynamics.

**Experimental setup.** We use proximal policy optimization (PPO) [13][2] to train RL policies due to its demonstrated effectiveness in robotic tasks and efficiency compared to other RL approaches. We consider four baselines that also assume unknown dynamics: certified control using contraction metric (C3M) [7], model-free PPO [13], SD-LQR [8], and standard LQR.

We consider four widely used simulation environments: 4D Car [7], 6D PVTOL [21, 17], 6D Neural-lander [22], and 10D Quadrotor [23, 17]. For real-world validation, we train policies in simulation using approximated dynamics of the TurtleBot3 Burger and deploy the best policy on a robot. Reference trajectories for real-world experiments were generated using a sinusoidal function, following the same procedure as in [7]. Detailed descriptions of our setup are provided in Appendix B.

---

[2]We follow the implementation details of [18] to ensure consistency across methods.

Table 1: Mean and 95% confidence intervals of MAUC (modified area under the curve) for normalized tracking error over 10 seeds. Mean inference time (ms) indicates per-step execution cost. The best result for each environment is bolded.

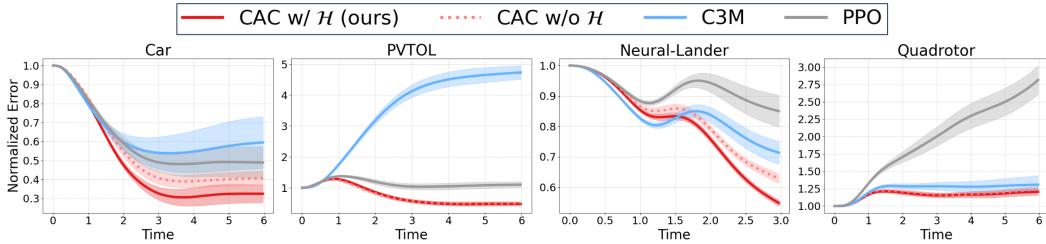| Method | MAUC (lower is better) | | | | Per Step Inference Time (ms) | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Car | PVTOL | NeuralLander | Quadrotor | Car | PVTOL | NeuralLander | Quadrotor |
| CAC (Ours) | $1.7 \pm 0.24$ | $\mathbf{2.7 \pm 0.38}$ | $\mathbf{2.0 \pm 0.05}$ | $\mathbf{7.1 \pm 0.41}$ | $\mathbf{0.07}$ | $\mathbf{0.08}$ | $\mathbf{0.09}$ | $\mathbf{0.08}$ |
| C3M | $2.3 \pm 0.56$ | $13 \pm 1.8$ | $2.2 \pm 0.16$ | $7.7 \pm 0.84$ | $0.1$ | $0.1$ | $0.1$ | $0.1$ |
| PPO [13] | $2.2 \pm 0.36$ | $4.4 \pm 0.57$ | $2.5 \pm 0.13$ | $14 \pm 1.8$ | $\mathbf{0.07}$ | $\mathbf{0.08}$ | $\mathbf{0.09}$ | $\mathbf{0.08}$ |
| SD-LQR [8] | $1.6 \pm 0.065$ | $13 \pm 1.7$ | $2.4 \pm 0.10$ | $7.5 \pm 0.35$ | $1$ | $10$ | $10$ | $20$ |
| LQR | $\mathbf{1.5 \pm 0.044}$ | $13 \pm 1.6$ | $2.7 \pm 0.17$ | $7.2 \pm 0.37$ | $4$ | $10$ | $10$ | $30$ |



Figure 2: We plot the the mean and 95% confidence intervals for normalized tracking error, $\|x(t) - x_d(t)\|_2 / \|x(0) - x_d(0)\|_2$, over time for each environment over 10 seeds.

**Performance metrics.** [7, 8] uses the area under the curve (AUC) of the normalized tracking error as a performance metric. However, to account for varying episode lengths, we adopt a *modified AUC (MAUC)* metric, defined as $MAUC = \frac{L}{T} \sum_{t=0}^{T} \frac{\|x(t) - x_d(t)\|_2}{\|x(0) - x_d(0)\|_2}$, where $L$ is the maximum episode length and $T$ is the length of the executed trajectory. After training, the best-performing policy (with the lowest MAUC) was selected for each seed. We report mean and 95% confidence intervals for MAUC evaluated over 10 training seeds, where each seed comprises 10 tracking trials for each of the 10 reference trajectories. We also report the mean per-step inference time for each method, measured on an Intel i7-1165G7 (2.80GHz) CPU with 16GB of 4266MHz LPDDR4X SDRAM.

**Performance evaluation.** Table 1 summarizes our simulation results, and Figure 2 presents normalized tracking errors over generated trajectories, comparing our method with and without entropy regularization—one of our key implementation details—along with key baselines. Additional results are given in Appendix D. We highlight four key observations:

*(1) Optimality of CAC:* For PVTOL, NeuralLander, and Quadrotor, CAC achieves the best performance, attaining the lowest MAUC among baseline algorithms. For the Car environment, the best performance is achieved by LQR, although CAC attains a comparably good result with only a slight margin. These results highlight the benefits of contraction-guided RL optimization in that it minimizes the *cumulative* tracking error. It is worthwhile to note that CAC consistently outperforms PPO, highlighting the benefit of integrating a contraction certificate into RL.

*(2) Robustness of CAC:* All baseline algorithms exhibit sensitivity to the underlying dynamics of the considered environments—for instance, C3M and LQR-based methods struggle with PVTOL, while PPO struggles with Quadrotor. In contrast, CAC maintains strong and consistent performance across all environments, potentially because our method uses model-free RL, which may inherently avoid issues related to model approximation errors—C3M and LQR heavily rely on the learned dynamics model. We further demonstrate this robustness in our real-world robot experiments (Section 5.1), where additional sim-to-real deviations are considered.

*(3) Practicality of CAC:* While the LQR-based methods show slightly lower MAUC than CAC for Car, they suffer from significant inference overhead due to solving a real-time optimization, resulting in inference times nearly $100\times$ longer than CAC. Moreover, inference time grows with problem dimensionality, indicating that these methods do not scale well. Such scaling is important, given
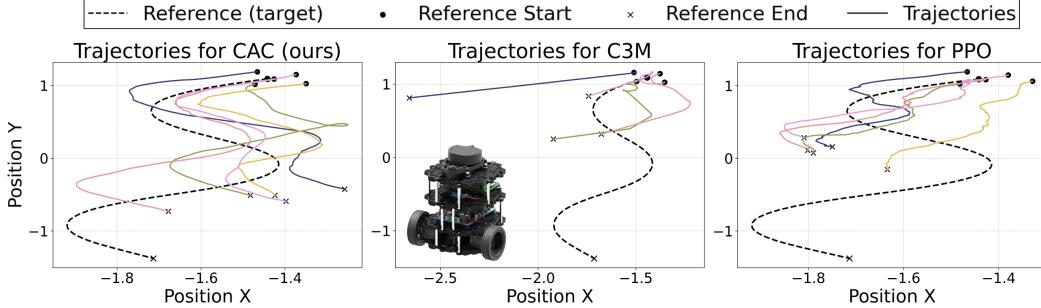
Figure 3: Robot demonstration of our method (CAC) compared with key baseline algorithms. The robot was initialized at five different initial locations and tasked to track the reference trajectory. The normalized tracking error across the horizon can be found in Appendix D.1.

that most robotic systems, such as drones and mobile robots, operate with limited computational capacity in high-dimensional environments.

*(4) Entropy-regularization of CAC:* As seen in Figure 2, removing entropy regularization of the CMG from CAC (dashed red) yields higher tracking error for Car and Neural-Lander, likely due to premature convergence to stationary points where the policy struggles to find an optimal contraction strategy.

## 5.1 Evaluation on Real-world Robots

We evaluate the sim-to-real transferability of our method by deploying policies (see Appendix B.8 for experiment details) trained in simulation on a TurtleBot3 Burger. Here, we consider both model approximation errors and real-world discrepancies, such as hardware imperfections and environmental uncertainties (e.g., traction coefficients). The path-tracking results are visualized in Figure 3 with 5 different initial positions, and the video is available in this link. We see that CAC outperforms C3M and PPO, where C3M completely fails to track the path and PPO exhibits significant divergence. In contrast, CAC maintains trajectory tracking within an acceptable margin, further demonstrating its robustness to model and real-world uncertainties that cause C3M and PPO to fail.

## 6 Related Works

Some prior works have explored learning control policies under contraction theory. [1] formulates contraction conditions for designing controllers as convex constraints using linear matrix inequalities (LMIs), where [5, 24] combine this with offline imitation learning to train contraction-preserving controllers. Other techniques use Sum-of-Squares (SoS) programming [6] or reproducing kernel Hilbert space (RKHS) theory [25] to synthesize CCMs, even with partially known dynamics. Despite their rigor, these approaches are limited by restrictive assumptions on system structure, focus solely on learning CCMs rather than minimizing cumulative tracking error, and require expensive geodesic computations to derive control inputs. SoS-based methods additionally demand polynomial dynamics and careful template tuning, which further hinders scalability.

## 7 Conclusions

We propose a contraction-certified RL algorithm, CAC, for robust path-tracking in settings with unknown dynamics. Our key idea is to simultaneously learn a CMG and control policy, where the CMG is used to sample a CCM that defines the reward function optimized by the policy. Our proposed method is scalable to complex systems with unknown dynamics, requires minimal design effort as it is end-to-end trainable, and, in principle, can be integrated with any RL algorithm. We theoretically show that our RL algorithm inherits a contraction certificate and extensively verify the benefits of our method through simulated and real-world experiments.

## Limitations

While our method demonstrates strong empirical performance, it has several limitations. First, the training process requires online interaction, which may be costly or impractical in certain real-world settings without reliable simulators. Second, although our theoretical results provide convergence guarantees under contraction conditions, satisfying these conditions in practice can be challenging. As a result, the validity of the theoretical guarantees may not always hold during training or execution.

## A    Theoretical Results

### A.1    Proof of Lemma 1

**Lemma** (Bound of Contraction Inequality). *If* (11) *holds, the performance measure in Equation* (12) *satisfies the following upper bound,*

$$\mathcal{J}_T^{\pi_c}(x_0, t_0) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \|\delta x(t_k)\|_M^2 \right] \le \frac{\|\delta x(t_0)\|_M^2}{1 - e^{2\lambda \Delta t}}.$$

*Proof.* Summing both sides of Equation (11) across the $\infty$ time steps in the trajectory under the contracting policy $\pi_c$ over the interval $[t_0, \infty]$, we obtain:

$$\mathcal{J}_T^\pi(x_0, t_0) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \|\delta x(t_k)\|_M^2 \right] \le \sum_{k=0}^{\infty} \|\delta x(t_0)\|_M^2 e^{-2\lambda k \Delta t}, \tag{15}$$

where $t_k = k\Delta t + t_0$ denotes the time at the $k^{th}$ sample.

We simplify the right-hand side by setting $C := \|\delta x(t_0)\|_M^2$ as,

$$\sum_{k=0}^{\infty} \|\delta x(t_0)\|_M^2 e^{-2\lambda k \Delta t} = C \sum_{k=0}^{\infty} e^{-2\lambda k \Delta t} = C \frac{1 - e^{-2\lambda \infty \Delta t}}{1 - e^{2\lambda \Delta t}} = \frac{C}{1 - e^{2\lambda \Delta t}}. \tag{16}$$

Therefore, we get,

$$\mathcal{J}_T^\pi(x_0, t_0) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \|\delta x(t_k)\|_M^2 \right] \le \frac{\|\delta x(t_0)\|_M^2}{1 - e^{2\lambda \Delta t}}. \tag{17}$$

$\square$

### A.2    Proof of Lemma 2

**Lemma** (Equivalence between Cost Minimization and Reward Maximization). *Define a cost function in terms of $\tilde{R}(x)$ as,*
$$C(x) = 1 - \tilde{R}(x) \in [0, 1].$$
*Then, the following equivalence holds,*

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \, \tilde{\mathcal{J}}_T^\pi(x_0, t_0) = \underset{\pi}{\operatorname{argmax}} \, \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k \tilde{R}_k \right] = \underset{\pi}{\operatorname{argmin}} \, \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k C_k \right].$$

*Proof.* The reward function defined in Equation (10) satisfies

$$\frac{1}{1 + \|\delta x\|_M^2} = 1 - \frac{\|\delta x\|_M^2}{1 + \|\delta x\|_M^2} \iff \tilde{R}(x) = 1 - C(x). \tag{18}$$

Therefore, the reward objective can be rewritten as

$$\mathcal{J}^\pi(x_0, t_0) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k \tilde{R}_k \right] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k (1 - C_k) \right]. \tag{19}$$

where $\tilde{R}(x(t_k))$ and $C(x(t_k))$ are abbreviated as $\tilde{R}_k$ and $C_k$, respectively.

Since the sum is finite due to the discount factor, we can separate the terms,

$$\mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k \left( 1 - C_k \right) \right] = \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k \right] - \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k C_k \right]. \tag{20}$$

Taking the maximization over a policy distribution yields

$$\begin{aligned}
\operatorname*{argmax}_\pi \mathcal{J}^\pi(x_0, t_0) &= \operatorname*{argmax}_\pi \left( \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k \right] - \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k C_k \right] \right) \\
&= \operatorname*{argmax}_\pi \left( \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k \right] \right) + \operatorname*{argmax}_\pi \left( -\mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k C_k \right] \right) \\
&= \operatorname*{argmax}_\pi \left( -\mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k C_k \right] \right) \\
&= \operatorname*{argmin}_\pi \left( \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k C_k \right] \right),
\end{aligned} \tag{21}$$

where we used that $\sum_{k=0}^\infty \gamma^k$ is independent of the policy $\pi$. $\qquad\square$

### A.3 Proof of Theorem 1

**Theorem** (Asymptotic Convergence of an RL-trained Policy via Existence of a Contracting Policy)**.**
*If there exists at least one contracting policy $\pi_c$ with a contraction rate $\alpha > 0$, then the optimal policy $\pi^*$ obtained via Equation (14) must exhibit asymptotic convergence, as given by,*

$$\lim_{t \to \infty} \| \delta x(t) \|_M^2 = 0, \ \forall x_0, t_0.$$

*Proof.* By Lemma 2, it suffices to prove cost minimization, as it is equivalent to reward maximization. Specifically, the cost term introduced in Lemma 2 admits the following upper bound,

$$\gamma \cdot C(x) = \gamma \frac{\| \delta x \|_M^2}{1 + \| \delta x \|_M^2} \leq C(x) \leq \| \delta x \|_M^2, \tag{22}$$

for any discount factor $\gamma \in (0, 1]$.

Denoting the time at the $k^{th}$ sample as $t_k = k\Delta t + t_0$ and the cost function at $t_k$ as $C(x(t_k)) := C_k$, we have the following inequality,

$$\mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k C_k \right] \leq \mathbb{E}_\pi \left[ \sum_{k=0}^\infty C_k \right] \leq \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \| \delta x(t_k) \|_M^2 \right] \leq \frac{\| \delta x(t_0) \|_M^2}{1 - e^{2\lambda \Delta t}}, \tag{23}$$

where the last inequality follows from Lemma 1.

Therefore, the infinite sum of the cost minimization objective is finite,

$$\mathbb{E}_\pi \left[ \sum_{k=0}^\infty \gamma^k C_k \right] \leq \mathbb{E}_\pi \left[ \sum_{k=0}^\infty C_k \right] \leq \mathbb{E}_\pi \left[ \sum_{k=0}^\infty \| \delta x(t_k) \|_M^2 \right] < \infty. \tag{24}$$

As a result, the cost minimization ensures asymptotic convergence by the Cauchy convergence theorem [26],

$$\lim_{t \to \infty} \| \delta x(t) \|_M^2 = 0. \tag{25}$$

$\qquad\square$

# B  Experimental Setup

We describe the dynamics model we used for the simulation in Section 5.

## B.1  Simulated Experiment Setup

We constructed a conventional RL environment to facilitate online training of both our method and the baseline algorithms, simulating each robot's dynamics using Euler's method. Each robot features distinct dynamics and predefined ranges for both state and control variables (i.e., a closed compact set). If the robot exceeds the state and control bounds, the episode is terminated and control inputs are clipped, respectively. Followed by CCM settings [1, 7], the algorithm is trained to generate the control correction $\delta u$, where the applied control is given by $u = u_d + \delta u$ with $u_d$ the reference trajectory's control, to facilitate the convergence of trajectories.

## B.2  Generation of Reference Trajectories

To evaluate tracking performance and facilitate RL training, we generate reference trajectories $x^*(t)$ using a randomized sinusoidal control synthesis method, similar to [7].

We assume a compact state space $\mathcal{X}$ for each environment and a subset $\mathcal{X}_0 \subset \mathcal{X}$ representing the initial state distribution. The initial reference state $x^*(0)$ is uniformly sampled from $\mathcal{X}_0$, and a corresponding control signal $u^*(t)$ is constructed as a linear combination of predefined sinusoidal signals:

$$u^*(t) = \sum_{i=1}^{n_f} w_i \sin\left(2\pi f_i \frac{t}{T}\right),$$

where $f_i$ are fixed frequencies, $w_i$ are weights sampled randomly per trajectory, and $T$ is the time horizon. The resulting control inputs are clipped to a fixed bound:

$$u^*(t) \in [0.75 \cdot u_{\min},\ 0.75 \cdot u_{\max}]$$

to ensure smoothness and prevent overly aggressive maneuvers, making the trajectory feasible for nearby tracking agents.

Given $x^*(0)$ and $u^*(t)$, the reference trajectory $x^*(t)$ is simulated forward using the system dynamics with zero disturbance (i.e., $d(t) = 0$) over a fixed horizon $t \in [0, T]$.

To define the initial condition for the closed-loop system, we sample an initial tracking error $x_e(0)$ uniformly from a compact error set $\mathcal{X}_e$. The actual initial state is then defined as:

$$x(0) = x^*(0) + x_e(0), \quad \text{where} \quad x^*(0) \in \mathcal{X}_0,\ x_e(0) \in \mathcal{X}_e.$$

## B.3  Car (4D)

We consider a control-affine car model by relaxing the fixed-velocity assumption of the classical Dubins car. The state is defined as $x := [p_x, p_y, \theta, v]$, where $p_x$ and $p_y$ denote the position, $\theta$ is the heading angle, and $v$ is the linear velocity. The control input is given by $u := [\omega, v]$, representing angular velocity and linear velocity, respectively.

The system dynamics follow a control-affine form:

$$\dot{x} = f(x) + B(x)u,$$

where

$$f(x) = \begin{bmatrix} v\cos(\theta) \\ v\sin(\theta) \\ 0 \\ 0 \end{bmatrix}, \quad B(x) = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

11

We define the following sets:

$$\mathcal{X} = \left\{ x \in \mathbb{R}^4 \,\middle|\, [-5, \ -5, \ -\pi, \ 1] \leq x \leq [5, \ 5, \ \pi, \ 2] \right\},$$
$$\mathcal{U} = \left\{ u \in \mathbb{R}^2 \,\middle|\, [-3, \ -3] \leq u \leq [3, \ 3] \right\},$$
$$\mathcal{X}_0 = \left\{ x_0 \in \mathbb{R}^4 \,\middle|\, [-2, \ -2, \ -1, \ 1.5] \leq x_0 \leq [2, \ 2, \ 1, \ 1.5] \right\},$$
$$\mathcal{X}_e = \left\{ x_{e0} \in \mathbb{R}^4 \,\middle|\, [-1, \ -1, \ -1, \ -1] \leq x_{e0} \leq [1, \ 1, \ 1, \ 1] \right\}.$$

## B.4   PVTOL (6D)

We refer the reader to the Appendix of [7] for a detailed description of the PVTOL environment.

## B.5   NeuralLander (6D)

We refer the reader to the Appendix of [7] for a detailed description of the NeuralLander environment.

## B.6   Quadrotor (10D)

We refer the reader to the Appendix of [7] for a detailed description of the Quadrotor environment.

## B.7   Simulated TurtleBot3 (3D)

The state of the TurtleBot3 system is defined as $x := [p_x, \ p_y, \ \theta]$, where $p_x$ and $p_y$ denote the planar position, and $\theta$ is the heading angle. The control input is $u := [\omega, \ v]$, representing the angular and linear velocities, respectively.

The system follows a control-affine dynamic model:

$$\dot{x} = f(x) + B(x)u,$$

where the drift term $f(x)$ is zero, and the input matrix $B(x)$ is defined as:

$$f(x) = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \quad B(x) = \begin{bmatrix} 0 & c_1 \cos(\theta) \\ 0 & c_2 \sin(\theta) \\ c_3 & 0 \end{bmatrix},$$

with empirical constants $c_1 = 0.9061$, $c_2 = 0.8831$, and $c_3 = 0.8548$, identified from real-world robot data using least squares regression. These coefficients capture the effects of friction and modeling uncertainties inherent to the TurtleBot3 Burger in our experimental setting.

The following sets are defined to reflect both the motion capture system of our facility and the operational specifications of the TurtleBot3 Burger:

$$\mathcal{X} = \left\{ x \in \mathbb{R}^3 \,\middle|\, [-5, \ -2, \ 0] \leq x \leq [0, \ 2, \ 2\pi] \right\},$$
$$\mathcal{U} = \left\{ u \in \mathbb{R}^2 \,\middle|\, [0, \ -1.82] \leq u \leq [0.22, \ 1.82] \right\},$$
$$\mathcal{X}_0 = \left\{ x_0 \in \mathbb{R}^3 \,\middle|\, [-1.7, \ 0.75, \ \pi] \leq x_0 \leq [-1.3, \ 1.15, \ 1.5\pi] \right\},$$
$$\mathcal{X}_e = \left\{ x_{e0} \in \mathbb{R}^3 \,\middle|\, [-0.1, \ -0.1, \ -0.25\pi] \leq x_{e0} \leq [0.1, \ 0.1, \ 0.25\pi] \right\}.$$

## B.8   Real-world TurtleBot3 (3D)

We employ a total of eight motion capture cameras to accurately track the position and orientation of the robot. Six of these cameras are visible in Figure 4, with two positioned behind the camera view. Our experiments begin by fixing the robot's initial position at the bottom-right corner of the map. A reference trajectory is then generated using a sinusoidal function, as described in Appendix B.2. For each trial, the robot is initialized with a slight offset from the trajectory's origin, and a pre-trained policy is deployed to track the reference path in real-time.
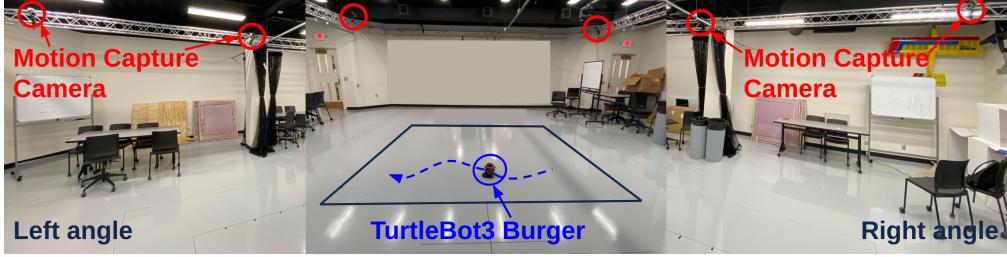
Figure 4: Our facility for real-world robot experiments. Qualisys motion capture cameras are circled with red color, and our robot is indicated in blue color circle.

## C   Algorithmic Implementation Details

Since our algorithms are trained in an online learning setting, we describe how to adapt offline C3M and LQR-based methods to operate in an online context.

### C.1   Dynamics Model Pre-training

In the unknown dynamics setting, CAC, C3M, SD-LQR, and LQR use a pre-trained dynamics model. Since our setting relies on an online simulator and is given the reference trajectory for every episode, we collect data by running reference controls with Gaussian noise. Then, the sample-wise MSE loss in Equation (8) for prediction of $(f, B)$ is minimized to $\sim 0.1$.

### C.2   Notes on Reward Function Design

It is worthwhile to note that reward design to directly satisfy the contraction condition in Equation (3) caused instability in learning a policy, which we attribute to the nonlinearity of the rewards and bias due to the approximated dynamics model. However, our proposed reward function, Equation (10), provides a principled direction for policy learning with strict concavity of the reward function while enforcing the contraction condition implicitly without explicit model dependence.

### C.3   C3M Implementation

The offline implementation of [7] is adapted to an online setting by changing C3M to train from online data collected after each update, similar to typical online training in RL. Note that we retain the structural assumptions about the environment present in their codebase[3] when adapting C3M to the online setting. These include the neural network architecture and the output design informed by environmental knowledge, even though we do not incorporate such assumptions ourselves.

### C.4   SD-LQR Implementation

SD-LQR leverages the state-dependent coefficient factorization (SDC) to reformulate the nonlinear drift dynamics in Equation (1) to the matrix multiplication as follows,

$$\dot{x} = f(x) + B(x)u = A(x)x + B(x)u, \tag{26}$$

which allows the controller to be expressed as $u = -R^{-1}B(x)^\top P(x)x$ and computed by solving the state-dependent Ricatti equation. That is, this formulation allows for the application of LQR to the time-varying path-tracking problem by formulating the error dynamics without linearization.

Since SD-LQR requires one to learn the factorization $A(x)$, we follow their implementation to obtain the required components using our dynamics model pre-training procedure. We refer the reader to [8] for more details.

---

[3] https://github.com/sundw2014/C3M

## C.5 Hyperparameters

We list the hyperparameters used to train CAC in Table 2.

**Table 2:** Summary of General, Dynamics, and CAC Hyperparameters

| Parameter | Value |
|---|---|
| *General Parameters* | |
| Optimizer | Adam [27] |
| CMG Upper Bound | 10.0 |
| CMG Lower Bound | 0.1 |
| *Dynamics Model Parameters* | |
| Network Layers | (256, 256) |
| Activation | Tanh |
| Batch Size | 1024 |
| *CAC Parameters* | |
| Network Layers (CMG, actor, critic) | (128, 128), (64, 64), (128, 128) |
| Learning Rate (CMG, actor, critic) | $1e{-}3$, $3e{-}4$, $1e{-}3$ |
| Activation Function (CMG, actor, critic) | Tanh, Tanh, Tanh |
| Batch size | $4{\times}256$ |
| Entropy Scaler (CMG, actor) | $1e{-}2$, $1e{-}2$ |
| K-Epochs | 10 |
| Target KL Divergence | 0.03 |
| GAE (Generalized Advantage Estimation) | 0.95 |
| Clipping Range | 0.2 |

# D  Additional Results

We present additional experimental results.

## D.1  Tracking Errors for Real-world Experiments

Appendix D.1 shows the normalized tracking error over a 20-second rollout of the TurtleBot3 Burger for key baselines.
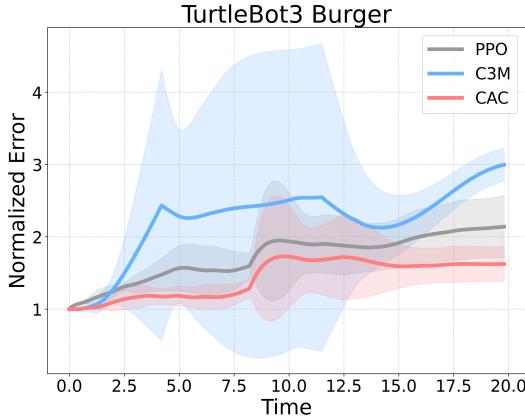


Figure 5: Mean and 95% confidence intervals for the normalized tracking error over 5 seeds from real-world experiments.

## D.2 Trajectory Visualizations for Simulated Experiments

We plot sampled trajectories for our method and key baselines across all environments in Figure 6. The model with the lowest MAUC among training seeds was used for each given reference trajectory.
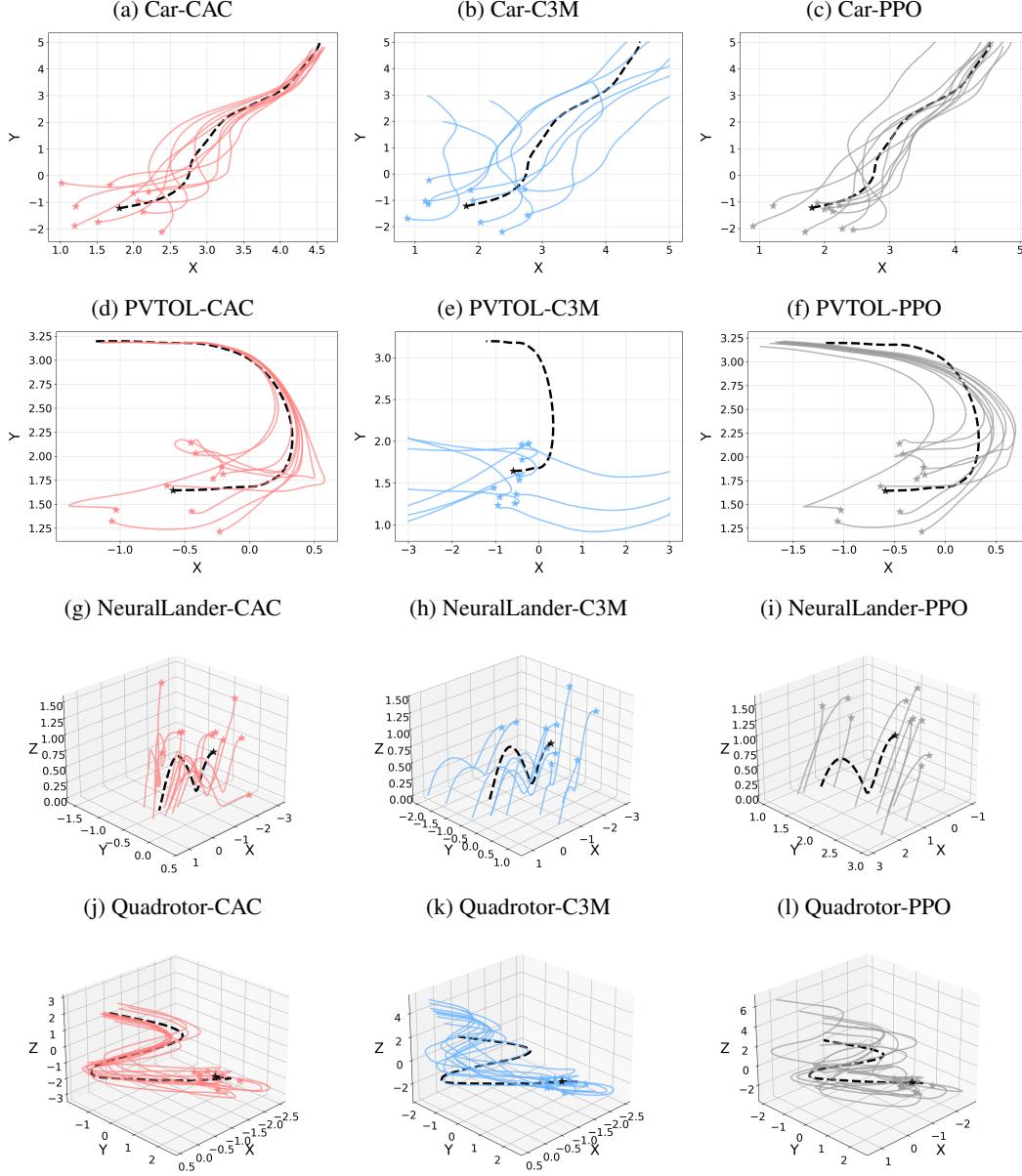


Figure 6: Path-tracking results for each algorithm with the same reference trajectory are depicted across environments.

# References

[1] I. R. Manchester and J.-J. E. Slotine. Control contraction metrics: Convex and intrinsic criteria for nonlinear feedback design. *IEEE Transactions on Automatic Control*, 62(6):3046–3053, 2017.

[2] W. Lohmiller and J.-J. E. Slotine. On contraction analysis for non-linear systems. *Automatica*, 34(6):683–696, 1998. ISSN 0005-1098. doi:https://doi.org/10.1016/S0005-1098(98)00019-3. URL https://www.sciencedirect.com/science/article/pii/S0005109898000193.

[3] A. Davydov, S. Jafarpour, and F. Bullo. Non-euclidean contraction theory for robust nonlinear stability. *IEEE Transactions on Automatic Control*, 67(12):6667–6681, 2022.

[4] H. Tsukamoto, S.-J. Chung, and J.-J. E. Slotine. Contraction theory for nonlinear stability analysis and learning-based control: A tutorial overview. *Annual Reviews in Control*, 52:135–169, 2021.

[5] H. Tsukamoto and S.-J. Chung. Neural contraction metrics for robust estimation and control: A convex optimization approach. *IEEE Control Systems Letters*, 5(1):211–216, 2020.

[6] S. Singh, B. Landry, A. Majumdar, J.-J. Slotine, and M. Pavone. Robust feedback motion planning via contraction theory. *The International Journal of Robotics Research*, 42(9):655–688, Aug. 2023. ISSN 0278-3649, 1741-3176. doi:10.1177/02783649231186165. URL https://journals.sagepub.com/doi/10.1177/02783649231186165.

[7] D. Sun, S. Jha, and C. Fan. Learning certified control using contraction metric. In *conference on Robot Learning*, pages 1519–1539. PMLR, 2021.

[8] S. M. Richards, J.-J. Slotine, N. Azizan, and M. Pavone. Learning control-oriented dynamical structure from data. In *International Conference on Machine Learning*, pages 29051–29062. PMLR, 2023.

[9] J.-J. E. Slotine. Modular stability tools for distributed computation and control. *International Journal of Adaptive Control and Signal Processing*, 17(6):397–416, 2003.

[10] M. L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.

[11] V. Konda and J. Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.

[12] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.

[13] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[14] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. Pmlr, 2018.

[15] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.

[16] Z. Ahmed, N. Le Roux, M. Norouzi, and D. Schuurmans. Understanding the impact of entropy on policy optimization. In *International conference on machine learning*, pages 151–160. PMLR, 2019.

[17] S. Singh, B. Landry, A. Majumdar, J.-J. Slotine, and M. Pavone. Robust feedback motion planning via contraction theory. *The International Journal of Robotics Research*, 42(9):655–688, 2023.

[18] S. Huang, R. F. J. Dossa, A. Raffin, A. Kanervisto, and W. Wang. The 37 implementation details of proximal policy optimization. In *ICLR Blog Track*, 2022. URL https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/. https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/.

[19] L. Engstrom, A. Ilyas, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, and A. Madry. Implementation matters in deep rl: A case study on ppo and trpo. In *International conference on learning representations*, 2019.

[20] M. Andrychowicz, A. Raichuk, P. Stańczyk, M. Orsini, S. Girgin, R. Marinier, L. Hussenot, M. Geist, O. Pietquin, M. Michalski, et al. What matters for on-policy deep actor-critic methods? a large-scale study. In *International conference on learning representations*, 2021.

[21] S. Singh, S. M. Richards, V. Sindhwani, J.-J. E. Slotine, and M. Pavone. Learning stabilizable nonlinear dynamics with contraction-based regularization. *The International Journal of Robotics Research*, 40(10-11):1123–1150, 2021.

[22] A. Liu, G. Shi, S.-J. Chung, A. Anandkumar, and Y. Yue. Robust regression for safe exploration in control. In *Learning for Dynamics and Control*, pages 608–619. PMLR, 2020.

[23] S. L. Herbert, M. Chen, S. Han, S. Bansal, J. F. Fisac, and C. J. Tomlin. Fastrack: A modular framework for fast and guaranteed safe motion planning. In *2017 IEEE 56th Annual Conference on Decision and Control (CDC)*, pages 1517–1522. IEEE, 2017.

[24] H. Tsukamoto and S.-J. Chung. Learning-based robust motion planning with guaranteed stability: A contraction theory approach. *IEEE Robotics and Automation Letters*, 6(4):6164–6171, 2021.

[25] S. Singh, V. Sindhwani, J.-J. E. Slotine, and M. Pavone. Learning Stabilizable Dynamical Systems via Control Contraction Metrics, Nov. 2018. URL http://arxiv.org/abs/1808.00113. arXiv:1808.00113 [cs].

[26] V. I. Bogachev and O. G. Smolyanov. *Real and functional analysis*. Springer, 2020.

[27] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.