# Critic as Lyapunov function (CALF): a model-free, stability-ensuring agent

Pavel Osinenko[1], Grigory Yaremenko[1], Roman Zashchitin[2],
Anton Bolychev[1], Sinan Ibrahim[1], Dmitrii Dobriborsci[2]

*Abstract*—**This work presents and showcases a novel reinforcement learning agent called Critic As Lyapunov Function (CALF) which is model-free and ensures online environment, in other words, dynamical system stabilization. Online means that in each learning episode, the said environment is stabilized. This, as demonstrated in a case study with a mobile robot simulator, greatly improves the overall learning performance. The base actor-critic scheme of CALF is analogous to SARSA. The latter did not show any success in reaching the target in our studies. However, a modified version thereof, called SARSA-m here, did succeed in some learning scenarios. Still, CALF greatly outperformed the said approach. CALF was also demonstrated to improve a nominal stabilizer provided to it. In summary, the presented agent may be considered a viable approach to fusing classical control with reinforcement learning. Its concurrent approaches are mostly either offline or model-based, like, for instance, those that fuse model-predictive control into the agent.**

## I. BACKGROUND AND PROBLEM STATEMENT

**Notation:** We will use Python-like array notation, e. g., $[0 : T] = \{0, \ldots, T-1\}$ or $s_{0:T} = \{s_0, \ldots, s_{T-1}\}$. Spaces of class kappa, kappa-infinity, kappa-ell and kappa-ell-infinity functions are denoted $\mathcal{K}, \mathcal{K}_\infty, \mathcal{K}\mathcal{L}, \mathcal{K}\mathcal{L}_\infty$, respectively. These are scalar monotonically increasing functions, zero at zero, and, additionally, tending to infinity in case of $\mathcal{K}_\infty$. The subscript $\geq 0$ in number set notation will indicate that only non-negative numbers are meant. The notation "cl $(\bullet)$", when referring to a set, will mean the closure. We denote modulo division of the first argument by the second as "$\bullet \bmod \bullet$".

### A. Problem statement

Consider the following optimal control problem:

$$V^\pi(s_0) = \sum_{t=0}^{\infty} c(s_t, a_t) \to \min,$$
$$s_0 \in \mathbb{S}, \pi \in \Pi, t \in \mathbb{Z}_{\geq 0} \tag{1}$$
$$\text{s.t.} \quad s_{t+1} = p(s_t, a_t), a_t \in \mathbb{A}(s_t)$$
$$a_t = \pi(s_t),$$

where:

[1]Skolkovo Institute of Science and Technology
[2]Deggendorf Institute of Technology, Technology Campus Cham
Corresponding author: P. Osinenko, email: p.osinenko@gmail.com.
First two authors contributed equally. Accepted for 27th IEEE Conference on Decision and Control.

1) $\mathbb{S}$ is the *state space*, e. g., $\mathbb{R}^n, n \in \mathbb{N}$, that is a normed vector space of all states of the given environment;
2) $\mathbb{A}(s)$ is the *action space*, in general state-dependent, e. g., a subset of $\mathbb{R}^m, m \in \mathbb{N}$, that is a set of all actions available to the agent at the current state. It is assumed that $\mathbb{A}(s)$ is compact for each state $s$;
3) $p : \mathbb{S} \times \mathbb{A} \to \mathbb{S}$ is the *state transition map* of the environment, assumed upper semi-continuous in norm;
4) $c : \mathbb{S} \times \mathbb{A} \to \mathbb{R}$ is the *cost function* of the problem, that is a function that takes a state $s_t$ and an action $a_t$ and returns the immediate cost $c_t$ incurred upon the agent if it were to perform action $a_t$ while in state $s_t$;
5) $\Pi$ is the *set of policies* (state feedback control laws in other words), that is a set of functions $\pi : \mathbb{S} \to \mathbb{A}$.

Sometimes, also a discount factor $\gamma \in (0, 1)$ is introduced in front of the cost in (1) as a factor $\gamma^t$.

The goal of reinforcement learning is to find a policy $\pi$ that minimizes $V^\pi(s_0)$ in (1). The policy $\pi^*$ that solves this problem is commonly referred to as *the optimal policy*. An *agent* will be referred to as a synonym for $\pi$ which was generated by some finite algorithm, e. g., actor-critic. In the latter, for instance, one step of the algorithm tries to approximate, i. e., "learn", the value $V^\pi$ via a neural network $\hat{V}^w$ with weights $w$ based on some collected data called *replay buffer*, say, $\mathcal{R} = \{s_{t_k}, a_{t_k}\}_{k \in [0:T_\mathcal{R}]}$ of size $T_\mathcal{R}$ and s. t. $t_k < t_{k+1}, \forall k$. In the second step, the algorithm may choose a so-called *greedy* or, e. g., $\varepsilon$-greedy action $a_t$, which minimize an actor loss $\mathcal{L}^{\text{act}}$, e. g., $\mathcal{L}^{\text{act}}(a_t) = c_t + \hat{V}^{w_t}(s_{t+1})$ or take a random action with probability $\varepsilon > 0$, respectively. A useful object is the action-value or Q-function $Q^\pi(s, a) \triangleq c(s, a) + V^\pi(s_+^a)$, where $s_+^a$ is the state at the next step relative to $s$ upon taking the action $a$. The herein present approach, CALF, will base particularly on Q-function.

Common reinforcement learning approaches to the above include various methods of policy gradient, actor-critics and tabular dynamic programming (see, e. g., [1]–[5]). A common classical control approach is in turn, e. g., model-predictive control (MPC) that places some finite horizon $T_{\text{pred}} > 2$ in place of the infinity in (1) and optimizes over respective sequences of actions via state propagation through the model $p$. A *model-free* agent, as will be understood in this work, is a controller that does not use $p$, or any learned model thereof, to compute actions. Such an agent is also called *data-driven* for instance. Notice, not every reinforcement learning agent is model-free, e. g., Dyna [6], [7].

Let now $\mathbb{G} \subset \mathbb{S}, 0 \in \mathbb{G}$ denote a *goal set* to which we want to drive the state optimally according to (1). We

assume, $\mathbb{G}$ to be a compact neighborhood of the origin. Also a compact in a subspace of the state space spanned by some state variables of interest may be considered, but we omit this case for simplicity. Let $d_{\mathbb{G}}(s) := \inf_{s' \in \mathbb{G}} \|s - s'\|$ be the distance to the goal. We call a policy $\pi_0 \in \Pi$ a $\mathbb{G}$-*stabilizer* if setting $a_t = \pi_0(s_t), \forall t$ implies that the distance between $s_t$ and $\mathbb{G}$ tends to zero over time. Formally, we use the following definition.

*Definition 1:* A policy $\pi_0 \in \Pi$ is called a $\mathbb{G}$-*stabilizer*, or simply a *stabilizer*, if the goal set is implied by the context, if

$$\forall t \geq 0 \ a_t \leftarrow \pi_0(s_t) \implies \forall s_0 \in \mathbb{S} \ \lim_{t \to \infty} d_{\mathbb{G}}(s_t) = 0. \quad (2)$$

It is called a uniform $\mathbb{G}$-stabilizer if, additionally, the limit in (2) is compact-convergent w.r.t. $s_0$ and

$$\exists \varepsilon_0 \geq 0 \ \forall t \geq 0 \ a_t \leftarrow \pi_0(s_t) \implies$$
$$\forall \varepsilon \geq \varepsilon_0 \ \exists \delta \geq 0 \ d_{\mathbb{G}}(s_0) \leq \delta \quad (3)$$
$$\implies \forall t \geq 0 \ d_{\mathbb{G}}(s_t) \leq \varepsilon,$$

where $\delta$ is unbounded as a function of $\varepsilon$ and $\delta = 0 \iff \varepsilon = \varepsilon_0$. The presence of an $\varepsilon_0$ means we do not insist on $\mathbb{G}$ being invariant. Thus, this extra condition only means that the state overshoot may be uniformly bounded over compacts in $\mathbb{S}$.

Notice that if the cost $c$ be continuous, strictly positive outside $\mathbb{G}$ and diverging with the distance to $\mathbb{G}$, the optimal policy $\pi^*$ is also a uniform $\mathbb{G}$-stabilizer.

The main hypothesis of this work is that incorporating some $\pi_0$ into the agent may improve its learning and guarantee stabilization into the goal set in all learning episodes, i.e., online. The main problem is: how to do this incorporation? Notice that even if $\pi_1, \pi_2$ are two stabilizers, an arbitrary switching or "shaping" onto into another will not provide stabilization guarantee in general.

## II. RELATED WORK

There are three principal methodologies for stability-ensuring reinforcement learning: shield-based reinforcement learning, the integration of MPC with reinforcement learning, and Lyapunov-based reinforcement learning.

Shield-based approaches involve a protective filter, also referred to as a shield, overseer, or supervisor, designed to prevent undesirable, say, destabilizing actions. These methods vary in how they evaluate actions and generate safe alternatives. Shields range from manual human oversight [8] to sophisticated formal verification variants [9]–[12]. Unlike human operators, formal logic shields are theoretically error-proof, but they require highly specific application development and can be complex, as detailed in [13]. Such techniques have been applied in areas such as probabilistic shielding [13], [14], supervisory systems for autonomous vehicles [15], and safe robotic manipulation [16]. However, human overseers introduce subjective biases and potential errors, and formal logic shields are often difficult to design.

The combination of MPC and reinforcement learning represents an active frontier in the quest for ensuring stability [17]–[22]. Such a fusion takes various forms, sometimes emphasizing model learning and other times focusing on safety constraints [23], [23]–[31]. The reinforcement learning

dreamer is a notable example that adopts the predictive spirit of MPC [32], [33]. Such approaches draw on the MPC's well-established ability to ensure closed-loop stability via techniques like terminal costs and constraints. Proposals such as the one by Zanon et al. [17], [18] embed robust MPC within reinforcement learning to maintain safety. Other research emphasizes the predictive control side more heavily [19], [20], building from a safe starting point towards a predictive model while maintaining the option to revert to safe policies when needed.

Lyapunov stability theory is well recognized in reinforcement learning, having roots dating back to the work of Perkins and Barto [34], [35], and has seen significant development since [20], [36]–[39]. Typically, these approaches are offline and require validation of a Lyapunov decay condition in the state space. Chow et al. [36], for instance, developed a safe Bellman operator to ensure Lyapunov compliance, while Berkenkamp et al. [20] used state space segmentation to validate Lyapunov conditions, demanding certain confidence levels in the statistical environment model. Online Lyapunov-based approaches also exist, often inspired by adaptive control techniques [40]. Robustifying controls, a key component introduced by [41], may distort the learning unless certain preconditions are met. Some reviews may be found in [42], [43]. Control barrier functions, another safety feature, have been successfully integrated with reinforcement learning, providing enhanced safety capabilities as seen in simulations with bipedal robots [44] and in model-free reinforcement learning agents [45]. Stochastic stability theory provides a basis for correctness [46], but the current landscape of Lyapunov-based methods often lacks capacity for real-time application without extensive pre-training, and is generally predicated on specific assumptions about environmental dynamics, such as second-order differentiability [47], linearity [48], or global Lipschitz continuity [49].

It should be stressed that policy shaping algorithms, e.g., [50] may be attractive in their similarity to pre-training the agent to boost learning, no online stabilization can be ensured under them. In contrast to the existing approaches, CALF is online, bringing the interplay of a stabilizer and the agent onto a rigorous footing, thus providing a viable, model-free approach to combining classical control and reinforcement learning. This is the essence of the current contribution.

## III. APPROACH

In general, it holds, due to the Hamilton-Jacobi-Bellman equation, that $\forall t \geq 0 \ Q_{t+1}^* - Q_t^* = -c_t^*$, where $Q_t^* = c(s_t, \pi^*(s_t)) + V^{\pi^*}\left(s_{t+1}^{\pi^*(s_t)}\right), c_t^* = c(s_t, \pi^*(s_t))$. This effectively means, given the conditions on the cost stated in Section I, that $Q_t^*$ is a Lyapunov function for the closed-loop system $s_{t+1} = p(s_t, \pi^*(s_t))$. If a model (critic) is employed, e.g., a deep neural network $\hat{Q}^w(s, a)$, due to imperfections of learning, the Lyapunov property may be lost, although desired. In CALF, we would like to retain the Lyapunov property of $\hat{Q}^w$. Enforcing Lyapunov-like constraints alone on $\hat{Q}^w$ would not solve the problem of ensuring stability, since those constraints may fail to be feasible at some time steps. Offline

approaches, as mentioned in Section II, would overcome this by large samples of state-action trajectories. What we do here instead is that we employ any stabilizer, let us call it $\pi_0$. The latter may be synthesized by common control techniques, like PID, sliding mode or funnel control. The question is, as stated earlier, how to combine the agent with $\pi_0$? This is brought into a systemic way in CALF. Namely, the critic update, i.e., update of the weights $w$ is done so as to try to satisfy Lyapunov-like constraints (decay and $\mathcal{K}_\infty$-bounds, see Algorithm 1, line 7). If this succeeds, the weights are passed and the next actor update will base upon them. If not, the weights are "frozen" (the respective variable is denoted $w^\dagger$) and $\pi_0$ is invoked. Along $w^\dagger$, the state-action pair is also "stored", namely, $s^\dagger, a^\dagger$. The Lyapunov-like constraints are checked relative to $s^\dagger, a^\dagger, w^\dagger$. This trick enables to safely combine the agent with $\pi_0$ so as to retain the stabilization guarantee of the latter. We use the Q-function to make the overall approach model-free. The actor loss $\mathcal{L}^{\text{act}}$ may be taken equal the current critic with last "successful" weights $w^\dagger$, namely, $\mathcal{L}^{\text{act}}(a) = \hat{Q}^{w^\dagger}(s_t, a)$. Any augmentation of the loss, e.g., penalties or entropy terms, may be included into $\mathcal{L}^{\text{act}}$, there is no restriction. Furthermore, one may choose to take an $\varepsilon$-greedy action.

Regarding the $\mathcal{K}_\infty$-bounds (see Algorithm 1, line 7), a reasonable choice of $\hat{\kappa}_{\text{low}}$, $\hat{\kappa}_{\text{up}}$ would be

$$\hat{\kappa}_{\text{low}}(\bullet) = C_{\text{low}}\bullet^2, \ \hat{\kappa}_{\text{up}}(\bullet) = C_{\text{up}}\bullet^2, \ 0 < C_{\text{low}} < C_{\text{up}}. \quad (4)$$

For the decay constraint, one may also take a quadratic rate or simply a constant $\bar{\nu} > 0$. Overall, the hyper-parameters $\hat{\kappa}_{\text{low}}$, $\hat{\kappa}_{\text{up}}$ and $\bar{\nu}$ determine a trade-off between freedom of learning and worst-case-scenario reaching time of the goal. If $\frac{\hat{\kappa}_{\text{low}}}{\hat{\kappa}_{\text{up}}}$ and $\bar{\nu}$ are chosen to be sufficiently small, the weights of the critic will not be prevented from converging to their ideal values, however if the critic is underfitted, smaller values of $\frac{\hat{\kappa}_{\text{low}}}{\hat{\kappa}_{\text{up}}}$, $\bar{\nu}$ may entail slower stabilization accordingly.

Now, the critic loss $\mathcal{L}^{\text{crit}}$ (see Algorithm 1, line 7) may be taken in various forms. The presented approach does not restrict the user. For instance, one may take a TD(1) on-policy loss as per:

$$\mathcal{L}^{\text{crit}}(w) = \sum_{k=0}^{T_{\mathcal{R}}} (\hat{Q}^w(s_{t_k}, a_{t_k}) - c(s_{t_k}, a_{t_k}) - \\ \hat{Q}^{w^\dagger}(s_{t_{k+1}}, a_{t_{k+1}}))^2 + \alpha_{\text{crit}}^{-2} \|w - w^\dagger\|^2. \quad (5)$$

The regularization term $\alpha_{\text{crit}}^{-2} \|w - w^\dagger\|^2$ is redundant if gradient descent based minimization is used, since one could simply set a learning rate $\alpha_{\text{crit}}$ as opposed to penalizing the displacement of weights. Notice the choice of the critic loss (or learning rate) does not prevent environment stabilization, although the quality of the learning may be affected. Finally, (5) is akin to the critic loss of SARSA due to its on-policy character, but this is not necessary, an off-policy loss may be used as well, e.g., with greedy actions instead of $a_{t_{k+1}}$ in (5).

*Remark 1:* An actor model $\pi^\theta$, e.g., as probability distribution, with weights $\theta$ may be employed instead of direct actions $a_t$ in Algorithm 1.

---

**Algorithm 1:** Critic as Lyapunov function (CALF) algorithm, model-free, action-value based

---

1: **Input**: $\bar{\nu} > 0, \hat{\kappa}_{\text{low}}, \hat{\kappa}_{\text{up}}, \pi_0$ : (uniform) stabilizer
2: **Initialize**: $s_0, a_0 := \pi_0(s_0), w_0$ s.t.

$$\hat{\kappa}_{\text{low}}(\|s_0\|) \leq \hat{Q}^{w_0}(s_0, a_0) \leq \hat{\kappa}_{\text{up}}(\|s_0\|)$$

3: $w^\dagger \leftarrow w_0, s^\dagger \leftarrow s_0, a^\dagger \leftarrow \pi_0(s_0), a_0 \leftarrow \pi_0(s_0)$
4: **for** $t := 1, \dots \infty$ **do**
5:      Take action $a_{t-1}$, get state $s_t$
6:      Update action: $a^* \leftarrow \arg\min_{a \in \mathbb{A}(s_t)} \hat{Q}^{w^\dagger}(s_t, a)$
7:      Try critic update

$$w^* \leftarrow \arg\min_{w \in \mathbb{W}} \mathcal{L}^{\text{crit}}(w) \\ \text{s.t. } \hat{Q}^w(s_t, a_t) - \hat{Q}^{w^\dagger}(s^\dagger, a^\dagger) \leq -\bar{\nu}, \\ \hat{\kappa}_{\text{low}}(\|s_t\|) \leq \hat{Q}^w(s_t, a_t) \leq \hat{\kappa}_{\text{up}}(\|s_t\|)$$

8:      **if** solution $w^*$ found **then**
9:          $s^\dagger \leftarrow s_t, a^\dagger \leftarrow a^*, w^\dagger \leftarrow w^*$
10:      **else**
11:          $a_t \leftarrow \pi_0(s_t)$
12:      **end if**
13: **end for**

---

### A. Modified SARSA

The new CALF agent was benchmarked via its immediate reinforcement learning alternative, *State–action–reward–state–action* (SARSA), which is essentially like Algorithm 1 prescribes (with on-policy critic loss), but with the Lyapunov-like constraints and the $\pi_0$ removed. In our case studies with a mobile robot, we observed such a plain SARSA failed to drive the robot to the target area within any reasonable number of learning iterations. To help SARSA succeed, we slightly modified it, namely, we retained the $w^\dagger$-mechanism, i.e., we used $w^\dagger$ in the critic loss (see Algorithm 1, line 7). We did not *enforce* the Lyapunov-like constraints in the optimization though. We only checked, whether those constraints were satisfied post factum and updated $w^\dagger$ accordingly as in CALF. Such a modification turned out to help SARSA reach the target in some learning runs. This algorithm will further be referred to as SARSA-m.

### IV. ANALYSIS

The main environment stability result is formulated in Theorem 1.

*Theorem 1:* Consider the problem (1) and Algorithm 1. Let $\pi_t$ denote the policy generated by Algorithm 1. If the policy $\pi_0$ is a stabilizer, then $\pi_t$ is a stabilizer. If the former is a uniform stabilizer, then so is $\pi_t$.

**Proof.**

First, let us consider $\mathbb{G}'$ to be a closed superset of $\mathbb{G}$, where the Hausdorff distance between $\mathbb{G}$ and $\mathbb{G}'$ is non-zero. Let $h$ denote the said distance.

Recalling Algorithm 1, let us denote:

$$\hat{Q}^\dagger := \hat{Q}^{w^\dagger}(s^\dagger, a^\dagger). \tag{6}$$

Next, we introduce:

$$\hat{\mathbb{T}} := \{t \in \mathbb{Z}_{\geq 0} : \text{successful critic update}\},$$
$$\mathbb{T}^{\pi_0} := \{0\} \cup$$
$$\left\{ t \in \mathbb{N} : \begin{array}{l} \text{successful critic update at } t-1 \wedge \\ \text{unsuccessful at } t \end{array} \right\}. \tag{7}$$

The former set represents the time steps at which the critic succeeds and the corresponding induced action will fire at the next step. The latter set consist of the time steps after each of which the critic first failed, discarding the subsequent failures if any occurred.

Now, let us define:

$$\hat{Q}_t^\dagger := \begin{cases} \hat{Q}_t, & t \in \hat{\mathbb{T}}, \\ \hat{Q}_{t-1}^\dagger, & \text{otherwise.} \end{cases}$$

Next, observe that there are at most

$$\hat{T} := \max\left\{ \frac{\hat{Q}_0^\dagger - \bar{\nu}}{\bar{\nu}}, 0 \right\} \tag{8}$$

steps until the critic stops succeeding and hence only $\pi_0$ is invoked from that moment on. Hence, $\hat{\mathbb{T}}$ is a finite set. Notice $\hat{T}$ was independent of $\mathbb{G}'$ and in turn dependent on the initial value of the critic.

Consider some $t^\dagger$, a time step after which the critic failed to find a solution. At step $t^\dagger + 1$, the action $a_{t^\dagger}$ is taken leading the state to transition into some $s_{t^\dagger + 1} = p(s_{t^\dagger}, a_{t^\dagger})$. Now, either the critic finds a solution $w_{t^\dagger + 1}^\dagger$ again, or $\pi_0$ is invoked.

In the latter case, by (2), let $T_{t^\dagger + 1}^{\pi_0}$ be s.t.

$$\forall t \geq t^\dagger + 1 + T_{t^\dagger + 1}^{\pi_0} \ d_{\mathbb{G}}(s_t) \leq h.$$

In other words, $\mathbb{G}'$ would be reached in no more than $\forall t \geq T_{t^\dagger + 1}^{\pi_0}$ steps. Notice $T_{t^\dagger + 1}^{\pi_0}$ implicitly depends on $s_{t^\dagger + 1}$.

Since there are at most $\hat{T}$ such episodes, where $\pi_0$ is invoked for at least one step, the set $\mathbb{T}^{\pi_0}$ is finite, although it implicitly depends on the initial state.

Let $\bar{T}^{\pi_0}$ be the maximal number of steps among $T_{t^\dagger + 1}^{\pi_0}, t^\dagger + 1 \in \mathbb{T}^{\pi_0}$ s.t.

$$\forall t \geq t^\dagger + 1 + T_{t^\dagger + 1}^{\pi_0} \ d_{\mathbb{G}}(s_t) \leq h.$$

That is, $\bar{T}^{\pi_0}$ is the longest time needed to reach $\mathbb{G}'$ starting from all, but finitely many states $s_{t^\dagger + 1}$ which are in turn uniquely determined by the initial state.

We thus conclude that the set $\mathbb{G}'$ is reached in no more than

$$T^* := \bar{T}^{\pi_0} \cdot \hat{T} \tag{9}$$

steps. Since $\mathbb{G}'$ was arbitrary, we conclude that $\pi_t$ is a stabilizer. Notice that the reaching time $T^*$ depends on the initial state $s_0$ and cannot, in general, be made uniform over compacts in $\mathbb{S}$ where the environment starts.

Now, let us address the case where $\pi_t$ can indeed be a uniform stabilizer. For this, we need to demonstrate uniformity of state overshoot and, respectively, uniformity of $T^*$ on compacts in states. To this end, let $\mathbb{S}_0$ be any compact in $\mathbb{S}$ and define

$$\bar{p}(\mathbb{S}_0) := \sup_{s \in \mathbb{S}_0, a \in \mathbb{A}(\mathbb{S}_0))} \|p(s, a)\|, \tag{10}$$

which exists since $p$ is upper semi-continuous in norm and $\mathbb{A}(s)$ is compact for every state $s$.

Define

$$s_t^\dagger := \begin{cases} s_t, & t \in \hat{\mathbb{T}}, \\ s_{t-1}^\dagger, & \text{otherwise.} \end{cases}$$

Observe that

$$\forall t \in \mathbb{Z}_{\geq 0} \ \left\| s_t^\dagger \right\| \leq \hat{\kappa}_{\text{low}}^{-1}(\hat{Q}_t^\dagger) \leq \hat{\kappa}_{\text{low}}^{-1}(\hat{Q}_0^\dagger).$$

Denote

$$\mathbb{S}_0^\dagger := \left\{ s \in \mathbb{S} : \|s\| \leq \sup_{s' \in \mathbb{S}_0} \hat{\kappa}_{\text{low}}^{-1}(\hat{\kappa}_{\text{up}}(\|s'\|)) \right\}. \tag{11}$$

Denote the state trajectory induced by the policy $\pi_0$ emanating from $s_0$ as $z_{0:\infty}^{\pi_0}(s_0)$ with single elements thereof denoted $z_t^{\pi_0}(s_0)$. By Proposition 2.2 in [51], since $\pi_0$ is a uniform stabilizer, there exists a $\mathcal{KL}$ function $\beta$ s.t.

$$\forall t \in \mathbb{Z}_{\geq 0}, s_0 \in \mathbb{S} \ d_{\mathbb{G}}(z_t^{\pi_0}) \leq \beta(d_{\mathbb{G}}(s_0), t). \tag{12}$$

By [52, Lemma 8], there exist two $\mathcal{K}_\infty$ functions $\kappa, \xi$ s.t.

$$\forall v > 0, t > 0 \ \beta(v, t) \leq \kappa(v)\xi(e^{-t}). \tag{13}$$

Since $0 \in \mathbb{G}$, it holds that

$$\forall s \in \mathbb{S} \ \kappa(d_{\mathbb{G}}(s)) \leq \kappa(\|s\|). \tag{14}$$

Hence

$$\forall s \in \mathbb{S}, t > 0 \ \beta(d_{\mathbb{G}}(s), t) \leq \kappa(\|s\|)\xi(e^{-t}). \tag{15}$$

It holds that

$$\forall t \in \mathbb{Z}_{\geq 0}$$
$$\beta(d_{\mathbb{G}}(s_t^\dagger), 0) \leq \kappa(d_{\mathbb{G}}(s_t^\dagger))\xi(1) \leq \kappa(\hat{\kappa}_{\text{low}}^{-1}(\hat{Q}_t^\dagger))\xi(1) \leq$$
$$\kappa(\hat{\kappa}_{\text{low}}^{-1}(\hat{Q}_0^\dagger))\xi(1)$$
$$\leq \kappa(\hat{\kappa}_{\text{low}}^{-1}(\hat{\kappa}_{\text{up}}(\|s_0\|)))\xi(1).$$

Let us define

$$\bar{\beta}^{\mathbb{S}_0^\dagger} := \sup_{\|s\| \leq \bar{p}(\mathbb{S}_0^\dagger)} \beta(s, 0).$$

If $\pi_0$ were invoked from $t^\dagger + 1$ on, the norm state would evolve bounded by $\beta(d_{\mathbb{G}}(s_{t^\dagger + 1}^\dagger), t - t^\dagger - 1), t \geq t^\dagger + 1$, due to (12), until either $\mathbb{G}'$ would be reached or the critic would succeed again. This together with a maximal possible "jump" of $\beta$ lets us to deduce

$$\forall t \ \beta(d_{\mathbb{G}}(s_{t^\dagger + 1}^\dagger), 0) \leq \beta(d_{\mathbb{G}}(s_t^\dagger), 0) + \bar{\beta}^{\mathbb{S}_0^\dagger}.$$

Now, let $\psi : \mathbb{R} \to \mathbb{R}$ be defined as:

$$\psi(v) := \kappa(\hat{\kappa}_{\text{low}}^{-1}(\hat{\kappa}_{\text{up}}(v)))\xi(1) + \bar{\beta}^{\mathbb{S}_0^\dagger} + \varepsilon_0. \tag{16}$$

We deduce that

$$\forall t \in \mathbb{Z}_{\geq 0} \ d_{\mathbb{G}}(s_t) \leq \psi(\|s_0\|).$$

Hence, we may claim the state resides in a desired $\varepsilon$-vicinity, with $\varepsilon \geq \varepsilon_0' := \bar{\beta}^{\mathbb{S}_0^\dagger} + \varepsilon_0$, around $\mathbb{G}$ if the initial state satisfies

$$\|s_0\| \leq \psi^{-1}(\varepsilon).$$

Notice the involved inverse exists due to the strict increase property.

Let $\varepsilon_{\mathbb{G}}$ be the Hausdorff distance between $\mathbb{G}$ and $\mathbb{G}'$. Now, we argue similarly to the non-uniform case as before, but defining

$$\bar{T}_{\text{unif}}^{\pi_0} := \sup_{\|s\| \leq \psi(\|s_0\|)}$$
$$\max\left\{1, -\log\left(\xi^{-1}\left(\frac{\varepsilon_{\mathbb{G}}}{\hat{\kappa}_{\text{low}}^{-1}(\hat{\kappa}_{\text{up}}(\|s\|))}\right)\right)\right\}.$$

Notice this definition is now uniform over a compact.

Notice also that it is straightforward to make $\hat{T}$ uniform over compacts in $\mathbb{S}$, e. g., by setting

$$\hat{T} := \max\left\{\frac{\sup_{\|s\| \leq \psi(\|s_0\|)} \hat{\kappa}_{\text{up}}(\|s\|) - \bar{\nu}}{\bar{\nu}}, 0\right\}.$$

The final reaching time reads $T^* = \hat{T} \cdot \bar{T}_{\text{unif}}^{\pi_0}$. ∎

*Remark 2:* If $\mathbb{G}$ is contained in the set $\mathbb{S}_{\bar{\nu}} := \{s \in \mathbb{S} : \hat{\kappa}_{\text{up}}(\|s\|) \leq \bar{\nu}\}$, then $\varepsilon_0' = \varepsilon_0$, in particular, it is zero if $\mathbb{G}$ is invariant. This claim is true since if the state is in $\mathbb{S}_{\bar{\nu}}$, then the critic will not fire and hence only $\pi_0$ is invoked leading to the same overshoot as under $\pi_0$. It is verified by setting

$$\psi(v) := \begin{cases} \kappa(\hat{\kappa}_{\text{low}}^{-1}(\hat{\kappa}_{\text{up}}(v)))\xi(1) + \bar{\beta}^{\mathbb{S}_0^\dagger} + \varepsilon_0, & v \geq \hat{\kappa}_{\text{up}}^{-1}(\bar{\nu}), \\ \kappa(v)\xi(1), & \text{otherwise.} \end{cases}$$

Now, set $\bar{\psi}(v) := \sup_{v' \leq v} \psi(v')$. It is easy to verify that $\bar{\psi}$ is an increasing function, hence it is Riemann integrable. Set

$$\Psi(v) := \int_{\varepsilon_0}^v \bar{\psi}(v') \, dv'.$$

It is easy to verify that $\Psi$ is continuous, strictly increasing and bounds $\bar{\psi}$ from above. Hence, we may claim that the state resides in a desired $\varepsilon$-vicinity, with $\varepsilon \geq \varepsilon_0$, around $\mathbb{G}$ if the initial state satisfies

$$\|s_0\| \leq \Psi^{-1}(\varepsilon).$$

The same argumentation works if one sets in line 8 of Algorithm 1 the condition as "**if** solution $w^*$ found and $s_t \notin \hat{\mathbb{G}}$" where $\hat{\mathbb{G}}$ is chosen to well contain $\mathbb{G}$. In that case, the critic deactivates in a vicinity of the goal. Notice otherwise Algorithm 1 does not rely on any explicit specification of the goal.

*Remark 3:* Since the number of invocations of $\pi_0$ is not greater than $T^*$ till the $\mathbb{G}$ is reached, the critic $\hat{Q}^\dagger$ is a multi-step Lyapunov function, i. e., $\hat{Q}^\dagger$ is non-increasing and

$$\forall t \in \mathbb{Z}_{\geq 0} \ \hat{Q}_{t+T^*}^\dagger - \hat{Q}_t^\dagger < 0.$$

*Remark 4:* For extensions of the stability result to stochastic environments and local basins of attraction, kindly refer to [53].

## V. CASE STUDY

### A. System description

We consider here a differential drive mobile wheeled robot (WMR) depicted in Fig. 1 for testing the proposed algorithms. The goal is to achieve autonomous stabilization of the robot at the origin of a coordinate system, starting from non-zero initial conditions. In the work space, there is a designated high-cost zone, which could represent a physical object like a swamp, a puddle, or a phenomenon that can adversely impact the robot's movement. Within this high-cost zone, the robot's maximum velocity is constrained to 1 cm per second to mitigate potential damage or difficulties in navigation.

The differential drive model describes the motion of the wheeled robot. In this model, the robot has two wheels mounted on either side, which can rotate independently. The position and orientation of the robot in the plane can be described by the state variables $(x, y, \vartheta)$, where:

- $x$ and $y$ represent the coordinates of the robot in the plane.
- $\vartheta$ is the orientation angle of the robot with respect to a reference direction.

The kinematic equations governing the motion of the robot are:

$$\begin{aligned} \dot{x} &= v\cos(\vartheta) \\ \dot{y} &= v\sin(\vartheta) \\ \dot{\vartheta} &= \omega \end{aligned} \tag{17}$$

where:

- $v$ is the linear velocity of the robot.
- $\omega$ is the angular velocity (rate of change of orientation).

In the high-cost zone, the linear velocity $v$ is restricted such that $v \leq 1$ cm/s. Outside of this zone, the robot can move at its maximum possible speed.

The control objective is to design a control law for $v$ and $\omega$ to stabilize the robot at the origin $(0, 0)$. This involves navigating from any initial position to the origin while possibly avoiding the high-cost zone. Experiments were conducted to validate the proposed algorithms under various initial conditions. The performance metrics typically include the time taken to reach the origin, the path taken by the robot, and its behavior in and around the high-cost zone.

The performance of CALF and SARSA-m as the benchmark agent were compared with a nominal stabilizer, and two MPC agents with different horizons. The latter are taken to get an idea about a nearly optimal controller for the stated problem. The code for the environment simulation, CALF, SARSA-m, MPC and nominal stabilizer may be found under https://github.com/osinenkop/rcognita-calf.

In our studies, we observed that a significant horizon of length 15 was sufficient to fulfill the stated goal. This controller is further referred to as MPC15. Notice that MPC15 may lack practicability due to such a long horizon, especially taking into account possible disturbances. Essentially, as the horizon length increases, the number of decision variables and constraints in the optimization problem grows, leading to increased computational complexity. Solving this problem within a reasonable time frame becomes more challenging as
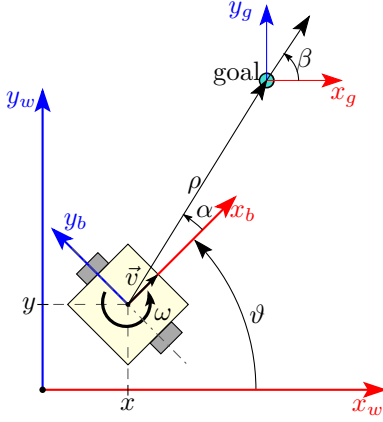
Fig. 1: Robot kinematics and its frames of interests.

the horizon grows. For real-time applications, such as controlling a differential drive mobile robot, the control decisions need to be computed and applied within very short time intervals. The extensive computations required for a 15-step horizon can exceed the available computational resources and time, making it possibly impractical for real-time implementation. In this regard, CALF and SARSA-m are beneficial as they do not use any models. Again, MPC15 was rather taken to get a clue of the best possible cost.

Recapping on the system description, at time $t$, the origin of the local coordinate frame (base) of the robot is located at the point $(x_t, y_t)$ of the world coordinate frame. The robot's linear velocity $v(t)$ along the axes determines the direction of motion, while the angular velocity $\omega(t)$ determines the rotation (refer to Fig. 1). The state vector $s \in \mathbb{R}^3$ and the control input vector (action) $a \in \mathbb{R}^2$ are defined as

$$s := \begin{bmatrix} x & y & \vartheta \end{bmatrix}^\top, a := \begin{bmatrix} v & \omega \end{bmatrix}^\top. \qquad (18)$$

The environment state transition map $p$ can be obtained via time discretization of the WMR differential equations which read: $\dot{x} = v\cos\vartheta, \dot{y} = v\sin\vartheta, \dot{\vartheta} = \omega$. For the studied MPC, we used Euler discretization, whereas for CALF and SARSA-m, $p$ is not necessary at all. In our case study, we used the controller sampling time of 100 ms.

The control actions were constrained by $a_{\min}$ and $a_{\max}$, which are determined by the actuators' capabilities. We took those constraints matching Robotis TurtleBot3, namely, normally (except for the high-cost spot) 0.22 m/s of maximal magnitude of translational and 2.84 rad/s of rotational velocity, respectively.

The control goal is to drive the robot into $\mathbb{G} = \{s \in \mathbb{R}^3 : \|s - s^*\| \leq \Delta\}$, where the target state $s^*$ was taken identical with the origin and $\Delta$ is target pose tolerance.

Regarding the high-cost zone, we introduced a spot on the plane $R^2$, namely, the cost was defined as (we did not penalize the action and only slightly penalized the angle):

$$c(s, a) = x^2 + y^2 + 0.1\vartheta^2 + 10c'(x, y), \qquad (19)$$

where

$$c'(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} \exp\left(-\frac{1}{2}\left(\frac{\Delta x^2}{\sigma_x^2} + \frac{\Delta y^2}{\sigma_y^2}\right)\right), \qquad (20)$$

with $\Delta x = x - \mu_x, \Delta_y = y - \mu_y$; $\sigma_x = \sigma_y = 0.1$ being standard deviations in the along-track and cross-track directions, respectively; and $\mu_x = -0.6, \mu_y = -0.5$ being the high cost spot center.

We now proceed to the description of the nominal stabilizer. The action components $v, \omega$ were determined based on the polar coordinate representation of the WMR as per [54], namely:

$$\begin{bmatrix} \dot{\rho} \\ \dot{\alpha} \\ \dot{\beta} \end{bmatrix} = \begin{bmatrix} \pm\cos\alpha & 0 \\ \mp\frac{\sin\alpha}{\rho} & 1 \\ \pm\frac{\sin\alpha}{\rho} & 0 \end{bmatrix} \begin{bmatrix} \pm v \\ \omega \end{bmatrix}, \qquad (21)$$

where the top sign (plus or minus) holds if $\alpha \in (-\frac{\pi}{2}, \frac{\pi}{2}]$ and the bottom sign holds if $\alpha \in (-\pi, -\pi/2] \cup (\pi/2, \pi]$. The transformation into the polar coordinates in turn reads:

$$\rho = \sqrt{x^2 + y^2},$$
$$\alpha = -\vartheta + \arctan 2(y, x),$$
$$\beta = -\vartheta - \alpha.$$

The stabilizer was set as per:

$$v \leftarrow K_\rho\rho, \qquad (22)$$
$$\omega \leftarrow K_\alpha\alpha + K_\beta\beta, \qquad (23)$$

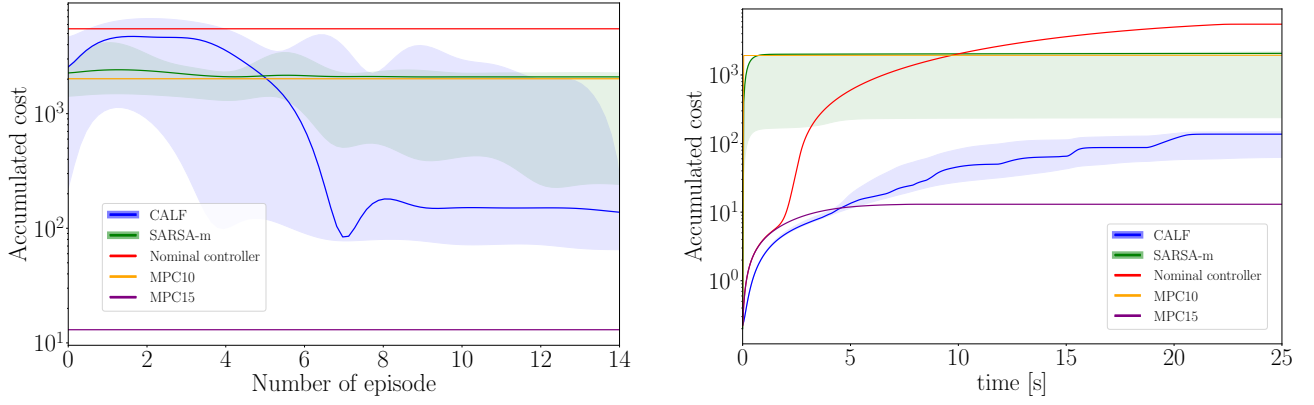where $K_\rho > 0, K_\beta < 0, K_\alpha - K_\rho > 0$.

The robot was run for 30 seconds in each episode. If the target was not reached, we nominally added a cost of 2000 to rule out the cases where the robot simply did not drive to the target while also not getting into the "hot" spot determined by $c'$. The target area was to 0.2 in radius around the origin (see Fig. 3). The initial robot pose was set to $\{-1 \quad -1 \quad \pi/2\}$.

### B. Discussion of results and conclusion

As can be seen from the results (Fig. 2, Fig. 3), both reinforcement learning agents outperformed MPC10 (CALF at once and SARSA-m in the final episode). The controller MPC15 performed best in terms of cost, but it should only be taken as a reference for a potential nearly best performing approach. It has a significant horizon length which reduces its practicability. Unlike MPC15, the studied reinforcement learning agents are totally model-free. Next, CALF also outperformed the nominal controller. What is remarkable is that both agents were able to detour the "hot" spot (notice in Fig. 3 how the nominal stabilizer is blind to the spot). CALF did it in all best learning episodes, whereas SARSA-m succeeded only in a part of those. It should be repeated here that plain SARSA did not succeed in our case studies at all. Overall, CALF always succeeded in reaching the target, thanks to its design and Theorem 1, and also outperformed the benchmark agent SARSA-m in terms of learning performance. This validates the initial claim that ensuring online environment stability is not just practicable, but it is also beneficial for episode-to-episode learning.

### REFERENCES

[1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction.* Cambridge, MA, USA: A Bradford Book, 2018.

(a) Learning curves depending on the episode, i. e., environment run number. The clouds represent the 95 % confidence intervals. The solid lines are median. Model-predictive controllers, as well as the nominal stabilizer, are given for reference.

(b) Accumulated cost (logarithmic scale) depending on time in the best learning episode. The clouds represent the 95 % confidence intervals. The solid lines are medians.

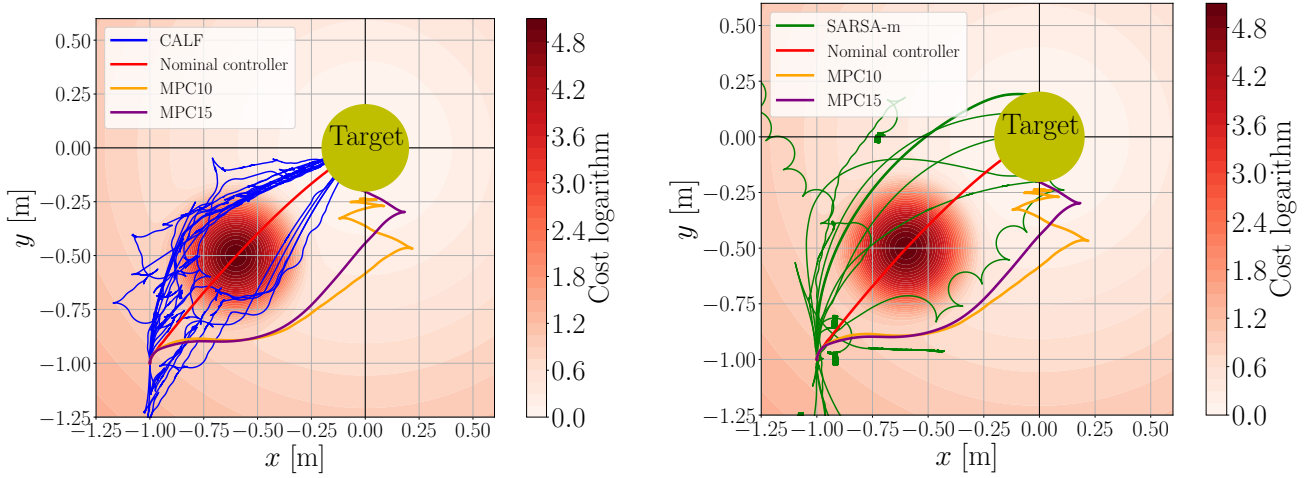Fig. 2: Learning curves obtained from 25 seeds of random number generator.



Fig. 3: The robot trajectories in best learning episodes over 25 seeds.

[2] S. M. Kakade, "A natural policy gradient," *Advances in neural information processing systems*, vol. 14, 2001.

[3] J. Baxter and P. L. Bartlett, "Infinite-horizon policy-gradient estimation," *Journal of Artificial Intelligence Research*, vol. 15, pp. 319–350, 2001.

[4] J. Peters and S. Schaal, "Policy gradient methods for robotics," in *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2006, pp. 2219–2225.

[5] D. P. Bertsekas, *Reinforcement learning and optimal control*. Athena Scientific Belmont, MA, 2019.

[6] R. S. Sutton, "Dyna, an integrated architecture for learning, planning, and reacting," *ACM Sigart Bulletin*, vol. 2, no. 4, pp. 160–163, 1991.

[7] M. Pei, H. An, B. Liu, and C. Wang, "An improved dyna-q algorithm for mobile robot path planning in unknown dynamic environment," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 52, no. 7, pp. 4415–4425, 2021.

[8] W. Saunders, G. Sastry, A. Stuhlmueller, and O. Evans, "Trial without error: Towards safe reinforcement learning via human intervention," *arXiv preprint arXiv:1707.05173*, 2017.

[9] Y. K. Tan and A. Platzer, "Deductive stability proofs for ordinary differential equations," *arXiv:2010.13096*, 2020.

[10] A. Platzer and J.-D. Quesel, "Keymaera: A hybrid theorem prover for hybrid systems (system description)," in *Automated Reasoning*. Springer, 2008, pp. 171–178.

[11] ——, "European train control system: A case study in formal verification," in *Formal Engineering Methods*. Springer, 2009, pp. 246–265.

[12] N. Fulton and A. Platzer, "Safe reinforcement learning via formal methods: Toward safe control through proof and learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[13] B. Könighofer, F. Lorber, N. Jansen, and R. Bloem, "Shield synthesis for reinforcement learning," in *International Symposium on Leveraging Applications of Formal Methods*. Springer, 2020, pp. 290–306.

[14] B. Könighofer, R. Bloem, S. Junges, N. Jansen, and A. Serban, "Safe reinforcement learning using probabilistic shields," in *International Conference on Concurrency Theory: 31st CONCUR 2020: Vienna, Austria (Virtual Conference)*. Schloss Dagstuhl-Leibniz-Zentrum fur Informatik GmbH, Dagstuhl Publishing, 2020.

[15] D. Isele, A. Nakhaei, and K. Fujimura, "Safe reinforcement learning on autonomous vehicles," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 1–6.

[16] B. Thananjeyan, A. Balakrishna, S. Nair, M. Luo, K. Srinivasan, M. Hwang, J. E. Gonzalez, J. Ibarz, C. Finn, and K. Goldberg, "Recovery rl: Safe reinforcement learning with learned recovery zones," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4915–4922, 2021.

[17] M. Zanon and S. Gros, "Safe reinforcement learning using robust MPC," *IEEE Transactions on Automatic Control*, vol. 66, no. 8, pp. 3638–3652, 2020.

[18] M. Zanon, S. Gros, and A. Bemporad, "Practical reinforcement learning of stabilizing economic mpc," in *2019 18th European Control Conference (ECC)*, 2019, pp. 2258–2263.

[19] T. Koller, F. Berkenkamp, M. Turchetta, and A. Krause, "Learning-based

model predictive control for safe exploration," in *2018 IEEE Conference on Decision and Control (CDC)*. IEEE, dec 2018.

[20] F. Berkenkamp, M. Turchetta, A. Schoellig, and A. Krause, "Safe model-based reinforcement learning with stability guarantees," in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds., vol. 30. Curran Associates, Inc., 2017.

[21] F. Berkenkamp, "Safe exploration in reinforcement learning: Theory and applications in robotics," Ph.D. dissertation, ETH Zurich, 2019.

[22] T. H. Oh, "Q-mpc: stable and efficient reinforcement learning using model predictive control," *IFAC-PapersOnLine*, vol. 56, no. 2, pp. 2727–2732, 2023.

[23] N. Karnchanachari, M. de la Iglesia Valls, D. Hoeller, and M. Hutter, "Practical reinforcement learning for mpc: Learning from sparse objectives in under an hour on a real robot," in *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, ser. Proceedings of Machine Learning Research, A. M. Bayen, A. Jadbabaie, G. Pappas, P. A. Parrilo, B. Recht, C. Tomlin, and M. Zeilinger, Eds., vol. 120. The Cloud: PMLR, 2020, pp. 211–224.

[24] K. Lowrey, A. Rajeswaran, S. Kakade, E. Todorov, and I. Mordatch, "Plan online, learn offline: Efficient learning and exploration via model-based control," *arXiv preprint arXiv:1811.01848*, 2018.

[25] W. Cai, A. B. Kordabad, and S. Gros, "Energy management in residential microgrid using model predictive control-based reinforcement learning and shapley value," *Engineering Applications of Artificial Intelligence*, vol. 119, p. 105793, 2023.

[26] B. Amos, I. Jimenez, J. Sacks, B. Boots, and J. Z. Kolter, "Differentiable MPC for end-to-end planning and control," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018, pp. 8289–8300.

[27] D. Hoeller, F. Farshidian, and M. Hutter, "Deep value model predictive control," in *Proceedings of the Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, L. P. Kaelbling, D. Kragic, and K. Sugiura, Eds., vol. 100, 2020, pp. 990–1004.

[28] S. East, M. Gallieri, J. Masci, J. Koutnik, and M. Cannon, "Infinite-horizon differentiable model predictive control," *International Conference on Learning Representations*, Jan. 2020.

[29] S. Reddy, A. D. Dragan, S. Levine, S. Legg, and J. Leike, "Learning human objectives by evaluating hypothetical behavior," *International Conference on Machine Learning*, 2019.

[30] C. Finn and S. Levine, "Deep visual foresight for planning robot motion," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 2786–2793.

[31] K. D. Asis, A. Chan, S. Pitis, R. Sutton, and D. Graves, "Fixed-horizon temporal difference methods for stable reinforcement learning," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, pp. 3741–3748, 2020.

[32] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi, "Dream to control: Learning behaviors by latent imagination," *International Conference on Learning Representations*, 2020.

[33] P. Wu, A. Escontrela, D. Hafner, P. Abbeel, and K. Goldberg, "Daydreamer: World models for physical robot learning," in *Conference on Robot Learning*. PMLR, 2023, pp. 2226–2240.

[34] T. Perkins and A. Barto, "Lyapunov design for safe reinforcement learning control," in *Safe Learning Agents: Papers from the 2002 AAAI Symposium*, 2002, pp. 23–30.

[35] T. J. Perkins and A. G. Barto, "Lyapunov-constrained action sets for reinforcement learning," in *ICML*, vol. 1, 2001, pp. 409–416.

[36] Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh, "A lyapunov-based approach to safe reinforcement learning," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018.

[37] A. B. Jeddi, N. L. Dehghani, and A. Shafieezadeh, "Memory-augmented lyapunov-based safe reinforcement learning: end-to-end safety under uncertainty," *IEEE Transactions on Artificial Intelligence*, 2023.

[38] M. Han, L. Zhang, J. Wang, and W. Pan, "Actor-critic reinforcement learning for control with stability guarantee," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6217–6224, 2020.

[39] Y.-C. Chang and S. Gao, "Stabilizing neural control using self-learned almost lyapunov critics," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 1803–1809.

[40] H. Zhang, L. Cui, X. Zhang, and Y. Luo, "Data-driven robust approximate optimal tracking control for unknown general nonlinear systems using adaptive dynamic programming method," *IEEE Transactions on Neural Networks*, vol. 22, no. 12, pp. 2226–2236, 2011.

[41] K. G. Vamvoudakis, M. F. Miranda, and J. P. Hespanha, "Asymptotically stable adaptive–optimal control algorithm with saturating actuators and relaxed persistence of excitation," *IEEE transactions on neural networks and learning systems*, vol. 27, no. 11, pp. 2386–2398, 2015.

[42] R. Kamalapurkar, P. Walters, J. Rosenfeld, and W. E. Dixon, *Reinforcement learning for optimal feedback control: A Lyapunov-based approach*. Springer, 2018.

[43] P. Osinenko, D. Dobriborsci, and W. Aumer, "Reinforcement learning with guarantees: a review," *IFAC-PapersOnLine*, vol. 55, no. 15, pp. 123–128, 2022, presented at IFAC Conference on Intelligent Control and Automation Sciences. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S2405896322010308

[44] J. Choi, F. Castaneda, C. J. Tomlin, and K. Sreenath, "Reinforcement learning for safety-critical control under model uncertainty, using control lyapunov functions and control barrier functions," *arXiv preprint arXiv:2004.07584*, 2020.

[45] R. Cheng, G. Orosz, R. M. Murray, and J. W. Burdick, "End-to-end safe reinforcement learning through barrier functions for safety-critical continuous control tasks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 3387–3395.

[46] R. Khasminskii and G. Milstein, *Stochastic Stability of Differential Equations*, ser. Stochastic Modelling and Applied Probability. Springer, 2011.

[47] S. Bhasin, R. Kamalapurkar, M. Johnson, K. G. Vamvoudakis, F. L. Lewis, and W. E. Dixon, "A novel actor-critic-identifier architecture for approximate optimal control of uncertain nonlinear systems," *Automatica*, vol. 49, no. 1, pp. 82–92, 2013.

[48] K. G. Vamvoudakis, "Q-learning for continuous-time linear systems: A model-free infinite horizon optimal control approach," *Syst. Control Lett.*, vol. 100, pp. 14–20, 2017.

[49] D. Vrabie, K. G. Vamvoudakis, and F. L. Lewis, *Optimal Adaptive Control and Differential Games by Reinforcement Learning Principles*. Institution of Engineering and Technology, 2012.

[50] H. Plisnier, D. Steckelmacher, J. Willems, B. Depraetere, and A. Nowé, "Transferring multiple policies to hotstart reinforcement learning in an air compressor management problem," *arXiv preprint arXiv:2301.12820*, 2023.

[51] Z.-P. Jiang and Y. Wang, "A converse lyapunov theorem for discrete-time systems with disturbances," *Systems & Control Letters*, vol. 45, no. 1, pp. 49–58, jan 2002.

[52] E. D. Sontag, "Comments on integral variants of ISS," *Systems and Control Letters*, vol. 34, no. 1-2, pp. 93–100, may 1998.

[53] P. Osinenko, G. Yaremenko, G. Malaniia, and A. Bolychev, "An actor-critic framework for online control with environment stability guarantee," *IEEE Access*, 2023. [Online]. Available: https://ieeexplore.ieee.org/document/10223230

[54] A. Astolfi, "Exponential stabilization of nonholonomic systems via discontinuous control," *IFAC Proceedings Volumes*, vol. 28, no. 14, pp. 661–666, 1995.