

Meta SAC-Lag: Towards Deployable Safe Reinforcement Learning via MetaGradient-based Hyperparameter Tuning

Homayoun Honari¹, Amir M. Soufi Enayati¹, Mehran Ghafarian Tamizi², Homayoun Najjaran^{1,2}

Abstract—Safe Reinforcement Learning (Safe RL) is one of the prevalently studied subcategories of trial-and-error-based methods with the intention to be deployed on real-world systems. In safe RL, the goal is to maximize reward performance while minimizing constraints, often achieved by setting bounds on constraint functions and utilizing the Lagrangian method. However, deploying Lagrangian-based safe RL in real-world scenarios is challenging due to the necessity of threshold fine-tuning, as imprecise adjustments may lead to suboptimal policy convergence. To mitigate this challenge, we propose a unified Lagrangian-based model-free architecture called *Meta Soft Actor-Critic Lagrangian* (Meta SAC-Lag). Meta SAC-Lag uses meta-gradient optimization to automatically update the safety-related hyperparameters. The proposed method is designed to address safe exploration and threshold adjustment with minimal hyperparameter tuning requirement. In our pipeline, the inner parameters are updated through the conventional formulation and the hyperparameters are adjusted using the meta-objectives which are defined based on the updated parameters. Our results show that the agent can reliably adjust the safety performance due to the relatively fast convergence rate of the safety threshold. We evaluate the performance of Meta SAC-Lag in five simulated environments against Lagrangian baselines, and the results demonstrate its capability to create synergy between parameters, yielding better or competitive results. Furthermore, we conduct a real-world experiment involving a robotic arm tasked with pouring coffee into a cup without spillage. Meta SAC-Lag is successfully trained to execute the task, while minimizing effort constraints. The success of Meta SAC-Lag in performing the experiment is intended to be a step toward practical deployment of safe RL algorithms to learn the control process of safety-critical real-world systems without explicit engineering.

I. INTRODUCTION

Reinforcement Learning (RL) is one of the most important paradigms for learning to control physical systems. However, a major shortcoming of RL is its need for exploration and extensive trial and error. For that reason, while we observe its wide success in various domains such as Energy systems [1], Video games [2], and Robotics [3], the real-world deployment of these algorithms to learn the control process poses is challenging [4] since the exploration process might lead the system to states that might damage the system and incur heavy costs to the user. To this end, safe RL methods aim to address this issue by optimizing the policy

*This work was supported by Apera AI and Mathematics of Information Technology and Complex Systems (MITACS) under IT16412 Mitacs Accelerate and Natural Sciences and Engineering Research Council (NSERC) Canada under the Grant RTI-2023-00418, and a partial equipment support was received from Kinova® Inc.

¹Department of Mechanical Engineering, University of Victoria, Victoria, BC, Canada {hmnhonari, amsoufi, najjaran}@uvic.ca

²Department of Electrical and Computer Engineering, University of Victoria, Victoria, BC, Canada mehrranght@uvic.ca

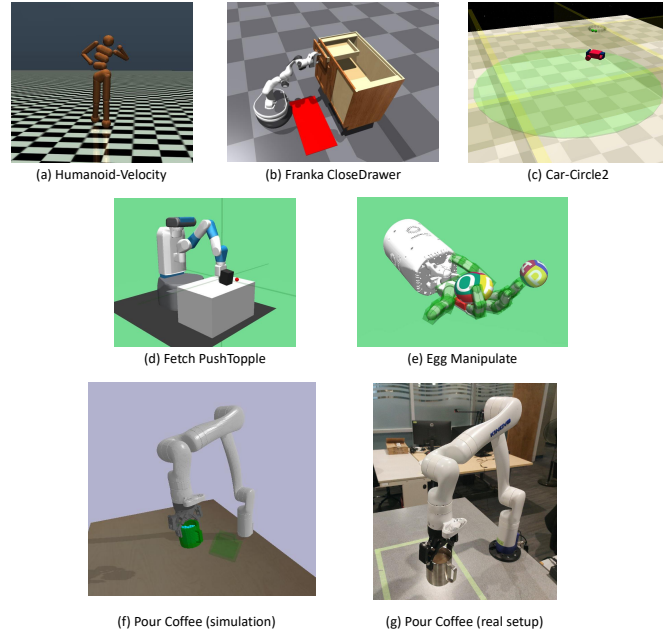


Fig. 1: Safety-critical environments used to deploy Meta SAC-Lag. The top two rows represent simulated environments with four general safety topics: locomotion (a), obstacle avoidance (b,c), robotic manipulation (d), dexterous manipulation (e). The bottom row represents Pour Coffee environment (f,g) used to study the deployability of the algorithm in a real-world setup.

such that it is compliant with the constraints. The constraints are defined such that they aim to prevent the system from exceeding its physical limitations.

The most common approach in safe RL is through the Lagrangian method. Specified under the Constrained Markov Decision Process (CMDP) framework, through defining thresholds for the constraints, the multi-objective optimization problem is converted to constraint satisfaction and is solved by casting it to an unconstrained problem using the Lagrangian method. While the approach has been extensively studied in the literature [5], [6], without precise tuning and engineering of the constraint thresholds, the Lagrangian methods will suffer from convergence to suboptimal policies. For that reason, the real-world use of these algorithms is rendered to be challenging due to the iterative process of hyperparameter tuning.

To address these challenges, based on the Soft Actor-Critic (SAC) architecture [7], our algorithm aims to address two fundamental problems: safe exploration and tuning-

free constraint adjustment. Previous attempts to automate the tuning of exploration-related hyperparameter of SAC have been mostly focused on optimizing the performance of the system [8], [9]. However, addressing the safety compliance of SAC has been limited. To this end, we propose a threshold-free safety-aware exploration optimization pipeline. In addition, our approach optimizes the safety threshold according to the overall performance of the policy. We are able to update the aforementioned hyperparameters using the metagradients w.r.t. to the meta-objectives which are computed based on the gradients of the internal learnable parameters. Finally, to assess the performance of Meta SAC-Lag, as depicted in Fig. 1, we study its performance against several baseline algorithms in five simulated robotic tasks with four different application themes. We observe that our method attains better or comparable results in terms of safety or reward performance while automatically tuning the safety-related hyperparameters. Furthermore, we present a safety benchmark test case, called *Pour Coffee*, which attempts to relocate and pour a coffee-filled mug into another cup. Constraint violation happens in case of collision or the spillage of the coffee. We deploy and train Meta SAC-Lag in a real-world setup using Kinova Gen3 robot. Our implementation shows that not only is Meta SAC-Lag capable of safe deployment without the iterative process of hyperparameter tuning but also, the learning process of the policy results in a smooth and jerk-free execution of the task with minimum effort imposed on the system.

Our main contributions can be summarized as:

- We propose a Lagrangian-based safe RL method able to automatically adjust the constraint bounds.
- Meta SAC-Lag addresses safe exploration through an unconstrained metagradient-based optimization pipeline.
- We validate the applicability of Meta SAC-Lag in five simulated robotic environments against baseline algorithms.
- A test environment, called *Pour Coffee*, is presented, and, with minimal prior safety-related hyperparameter tuning, Meta SAC-Lag is trained on a real-world Kinova Gen3 setup. The algorithm successfully achieves the task objective with minimized effort exerted on the robot.

II. RELATED WORK

The constrained Markov decision process (CMDP), is the theoretical building block of safe RL. CMDPs have been widely studied in the RL paradigm [10], [11] and are solved using Lagrangian methods [12]. In this regard, Shen et. al [13] devised the risk-sensitive policy optimization (RSPO) algorithm which sequentially decreases the Lagrangian multiplier to zero. Furthermore, Stooke et. al [14] updates the multiplier using PID control. Additionally, Reward-constrained policy optimization (RCPO) employs dual gradient descent optimization for the policy and Lagrange multiplier [15].

In another aspect, metagradient optimization has been explored thoroughly in RL hyperparameter tuning. Initially, model-agnostic meta-learning (MAML) [16] introduced meta-optimization of initial weights to enable fast

task adaptation within a few gradient descent steps. In a different approach, Meta-Gradient RL [17] extended the concept to learn the hyperparameters of return functions online. This paradigm offered a general approach, applicable to other RL hyperparameters. Subsequently, similar techniques were applied for auto-tuning other RL hyperparameters, such as exploration thresholds [18], entropy temperature in SAC [9], auxiliary tasks and sub-goals [19], and differentiable hyperparameters of loss functions [20]. Despite these advancements, metagradient methods have not been extensively explored in constrained RL paradigms [5], with few applications focusing on ensuring safety in sensitive learning tasks, as seen in the work by Calian et al. [21]. The authors utilized meta-gradients to update the Lagrange multiplier learning rate in an off-policy RL framework.

The Lagrangian methods are not the sole approach taken toward solving safe RL. Thananjeyan et al. [22] trained a recovery policy in parallel to the task policy and used it whenever the task policy chose actions deemed too risky. More prominently, model-based RL and safety guarantees for risk aversion during training are proposed. Koppejan et al. [23] used the neuroevolutionary approach to exploiting domain expertise to learn safe models for model-based RL. Thomas et al. [24], in a different approach, used near-future imagination to plan safe trajectories ahead of time. Moldovan et al. [25] focused on risk aversion in MDPs using near-optimal Chernoff bounds. Lyapunov functions have also been used to guarantee safety during training [26], [27], though constructing Lyapunov functions remains a challenge due to their typically hand-crafted nature and the absence of clear principles for agent safety and performance optimization.

In summary, RL agent evaluation relies on human-provided rewards, but the risk of misstated human objectives is often overlooked. Also, despite significant progress in safe RL methodologies, there remains a big gap in readily deployable agents in industrial contexts. Moreover, a critical balance exists between reward and cost in RL, as each action can impact both aspects, creating a multi-dimensional problem. These challenges hinder robust and dependable RL algorithms suitable for real-world implementation. Our motivation for this work is to use metagradient optimization for self-tuning of the safety threshold. This will minimize the need for safety-related hyperparameter tuning in safe RL while improving performance. Ultimately, the self-tuning safety threshold will enable us to deploy the agent directly in the real world.

III. BACKGROUND

In this section, we investigate the background by exploring essential preliminary concepts that serve as the foundation for this paper. We start with the discussion of the CMDP framework. Furthermore, we delve into the formulation of the safety critic and SAC.

A. Constrained Markov Decision Process (CMDP)

CMDP comprises the tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, c, \gamma_r, \gamma_c, \rho_0 \rangle$ where \mathcal{S} denotes the state space, \mathcal{A} represents the action

space, and r denotes the reward function: $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}$. The transition function $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto [0, 1]$ defines the likelihood $\mathcal{P}(s'|s, a)$ of moving from state s to s' by executing action a . The probability distribution function $\rho_0 : \mathcal{S} \mapsto [0, 1]$ denotes initial state distribution of the framework. Furthermore, $c(s)$ is the constraint indicator function which determines whether state s violates the constraint functions specified by C : $c(s) = \mathbb{1}[C(s) == 1]$. Parameters $\gamma_r \in [0, 1)$ and $\gamma_c \in [0, 1)$ serve as discount factors for reward and safety critics, respectively. Ultimately, the solution to CMDP is represented by the policy $\pi : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ which is the probability distribution over actions. The value function associated with policy π for a specific state-action pair (s, a) and the corresponding recursive equation, known as the Bellman equation, can be formulated as follows:

$$\begin{aligned} Q_r^\pi(s, a) &= \mathbb{E}_{s_t \sim \mathcal{P}, a_t \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right] \\ &= \mathbb{E}_{s' \sim \mathcal{P}} [r(s, a) + \gamma V_r^\pi(s')] \end{aligned} \quad (1)$$

Additionally, the primary function of the safety critic is to estimate the probability of a policy failure occurring in the future, determined by the expected cumulative discounted probability of failure.

$$\begin{aligned} Q_c^\pi(s, a) &= \mathbb{E}_{s_t \sim \mathcal{P}, a_t \sim \pi} \left[c(s) + (1 - c(s)) \sum_{t=1}^{\infty} [\gamma_c^t c(s_t)] \right] \\ &= \Pr[c(s) == 1] + \gamma_c \mathbb{E}_{s' \sim \mathcal{P}} [(1 - c(s)) V_c^\pi(s')] \end{aligned} \quad (2)$$

Finally, the main objective of an RL algorithm in a CMDP framework is to find a policy to maximize expected return while satisfying the constraints starting from the initial state s_0 :

$$\begin{aligned} \pi^* &= \operatorname{argmax}_{\pi \in \Pi} \mathcal{J}_r^\pi = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{s_0 \sim \rho_0} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \\ \text{s.t. } \mathcal{J}_c^\pi &= \mathbb{E}_{s_0 \sim \rho_0} \left[\sum_{t=0}^{\infty} \gamma_c^t c_t \right] \leq \varepsilon \end{aligned} \quad (3)$$

B. Soft Actor Critic (SAC)

SAC [7] optimizes a stochastic policy in an off-policy manner, utilizing two neural networks: one for estimating Q -function (critic) and another for policy updates (actor). A key feature of SAC is entropy regularization, where the policy aims to strike a balance between maximizing expected return and maximizing entropy. This balance mirrors the exploration-exploitation trade-off; higher entropy encourages greater exploration, potentially accelerating learning and prevent convergence to suboptimal solutions.

Considering ω_r and ϕ as parameters representing the critic and actor networks, respectively, training these networks involves sampling a batch of samples from the replay buffer. ω_r is updated by taking the gradient through the mean squared error (MSE) loss between the critic output and the target value:

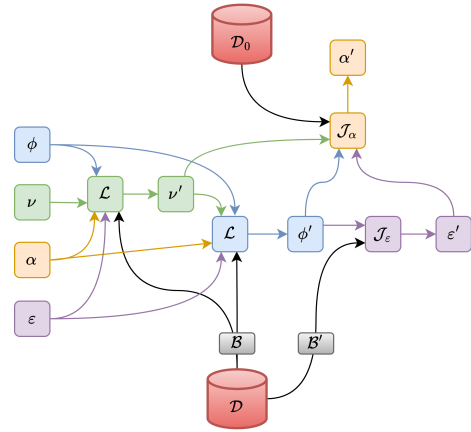


Fig. 2: The Computational Graph of the Meta SAC-Lag.

$$\mathcal{J}_r^{Q_{\omega_r}} = \mathbb{E}_{(s,a,r) \sim \mathcal{D}} \left[\frac{1}{2} (Q_{\omega_r}(s, a) - Q_r^{\text{tar}}(s, a))^2 \right], \quad (4)$$

where Q_r^{tar} is calculated as:

$$\begin{aligned} Q_r^{\text{tar}}(s, a) &= \\ &\mathbb{E}_{s' \sim \mathcal{P}(s,a), a' \sim \pi_\phi} [r(s, a) + \gamma_r (Q_r(s', a') - \alpha \log(\pi_\phi(a'|s')))] \end{aligned} \quad (5)$$

Furthermore, the policy π_ϕ is optimized by taking the gradient through the critic and the expected entropy of the policy:

$$\mathcal{J}_r^{\pi_\phi} = \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_\phi} [\alpha \log(\pi_\phi(a|s)) - Q_{\omega_r}(s, a)] \quad (6)$$

Finally, it is important to note that the safety critic (Q_{ω_c}) defined in Section III-A is trained using the same loss formulation in Eq. 4, without the entropy term.

IV. METHOD

In this section the process of metagradient optimization of the safety threshold ε and entropy temperature α is discussed.

A. Metagradient Optimization

Metagradient optimization is the process with which we can optimize the hyperparameters that are not a part of the main loss function. Fundamentally, these meta-parameters¹ dictate the dynamics of the system and direct it toward a certain behavior. In the context of metagradient reinforcement learning [17], in abstract terms, the learnable system variables are parameterized as θ . These parameters are updated to θ' by following the rule:

$$\theta' = \theta + f(\mathcal{J}, \theta, \eta, \mathcal{B}) \quad (7)$$

where η is the list of hyperparameters, \mathcal{B} a mini-batch of experience, and f the gradient of the objective function \mathcal{J} w.r.t. θ . Furthermore, the optimization process of the meta-parameters η can be formulated based on the updated parameter θ' :

$$\eta' = \eta + \beta_\eta \frac{\partial \mathcal{J}'(\theta', \eta, \mathcal{B}')}{\partial \eta} = \eta + \beta_\eta \frac{\partial \mathcal{J}'(\theta', \eta, \mathcal{B}')}{\partial \theta'} \frac{d\theta'}{d\eta} \quad (8)$$

¹In this paper, we use hyperparameters and meta-parameters terms interchangeably.

where \mathcal{J}' is the meta-objective used for the optimization of the meta-parameters, β_η the learning rate associated with η , and \mathcal{B}' a resampled mini-batch validation set similar to the cross-validation method in the meta-optimization literature [28]. Finally, $\frac{d\theta'}{d\eta}$ can be calculated as:

$$\frac{d\theta'}{d\eta} = \left(I + \frac{\partial f(\mathcal{J}, \theta, \eta, \mathcal{B})}{\partial \theta} \right) \frac{d\theta}{d\eta} + \frac{\partial f(\mathcal{J}, \theta, \eta, \mathcal{B})}{\partial \eta} \quad (9)$$

B. SAC-Lagrangian

In the Lagrangian version of the SAC we aim to optimize the policy based on its reward objective such that it is compliant with the safety objective:

$$\begin{aligned} \pi_\phi^* = \max_{\pi_\phi \in \Pi} \mathcal{J}_{r\pi_\phi}^{\pi_\phi} &= \mathbb{E}_{s \sim \mathcal{D}} [Q_{\omega_r}(s, a) - \alpha \log \pi_\phi(s, a)] \\ \text{s.t. } \mathcal{J}_c^{\pi_\phi} &= \mathbb{E}_{s \sim \mathcal{D}} [Q_{\omega_c}(s, a)] \leq \varepsilon \end{aligned} \quad (10)$$

Naturally, multiple constraints can be defined for the policy to consider all of them. However, in this paper, in order to keep the formulation simple and general, we consider a single constraint signal that is the result of the superposition of all the constraint functions. In this paper, in contrast to [7], we refrain from considering α as an additional constraint and aim to optimize it through metagradient optimization.

Furthermore, the optimization process of policy in Eq. 10 is formulated by casting it as a Lagrangian loss and back-propagating through the loss:

$$\begin{aligned} \min_{\nu \geq 0} \max_{\pi_\phi \in \Pi} \mathcal{L}(\pi_\phi, \nu, \varepsilon, \alpha) &= \mathcal{J}_{r\pi_\phi}^{\pi_\phi} - \nu(\mathcal{J}_c^{\pi_\phi} - \varepsilon) \\ &= \mathbb{E}_{s \sim \mathcal{D}} [Q_{\omega_r}(s, a) - \alpha \log \pi_\phi(s, a) - \nu(Q_{\omega_c}(s, a) - \varepsilon)] \end{aligned} \quad (11)$$

where ν is the Lagrange multiplier.

C. Meta SAC-Lag

Following the conventional notation in the context of gradient-based hyperparameter optimization [29], we split the parameters into *inner* and *outer* parameters. Rather than a one-shot optimization as in Eq. 7, we propose a sequential updating approach. We define and update the inner parameters as:

$$\theta'_{\text{inner}} = \begin{bmatrix} \nu' \\ \phi' \end{bmatrix} = \begin{bmatrix} \nu \\ \phi \end{bmatrix} + \begin{bmatrix} -\nabla_\nu \mathcal{L}(\pi_\phi, \nu, \varepsilon, \alpha) \\ \nabla_\phi \mathcal{L}(\pi_\phi, \nu, \varepsilon, \alpha) \end{bmatrix} \quad (12)$$

Furthermore, in the same sequential manner, we first update ε and then α :

$$\theta'_{\text{outer}} = \begin{bmatrix} \varepsilon' \\ \alpha' \end{bmatrix} = \begin{bmatrix} \varepsilon \\ \alpha \end{bmatrix} + \begin{bmatrix} \nabla_\varepsilon \mathcal{J}_\varepsilon(\pi_{\phi'}) \\ \nabla_\alpha \mathcal{J}_\alpha(\pi_{\phi'}, \nu', \varepsilon') \end{bmatrix} \quad (13)$$

where \mathcal{J}_ε and \mathcal{J}_α correspond to the objective functions of ε and α , respectively. To this end, we intended to design the objective function for ε solely based on the performance of the resultant policy. Our intuition behind the aforementioned design stems from the idea that the threshold should be adjusted such that it improves the performance of the agent as a whole. For that purpose, the ε objective function is proposed as:

$$\mathcal{J}_\varepsilon(\pi_{\phi'}) = \mathbb{E}_{s \sim \mathcal{D}} [\nu'_{\text{copy}} Q_{\omega_c}(s, a) - Q_{\omega_r}(s, a)] \quad (14)$$

Algorithm 1 Meta SAC-Lag

Require:

- Initialize Policy network ϕ^0 , Exploration rate α^0
 - Critic network $\omega_{r_1}^0, \omega_{r_2}^0$, Safety critic network $\omega_{c_1}^0, \omega_{c_2}^0$
 - Lagrangian values ε^0, ν^0
 - Learning rates $\beta_\phi, \beta_\varepsilon, \beta_\nu, \beta_\alpha$
 - 1: Create Transition buffer \mathcal{D} , Safety buffer \mathcal{D}_s , and Initial state buffer \mathcal{D}_0
 - 2: Randomly sample initial state $s_0 \sim \rho_0$ and fill \mathcal{D}_0
 - 3: **for** $e = 1, \dots$ **do**
 - 4: Reset environment $s_0 \sim \rho_0 = \text{env.reset}()$
 - 5: **for** $t = 0, \dots, T - 1$ **do**
 - 6: Sample action $a_t \sim \pi_\phi$
 - 7: $s_{t+1}, r_t, c_t \leftarrow \text{env.step}(a_t)$
 - 8: **if** $c_t == 1$ **then**
 - 9: $\mathcal{D}_s \leftarrow \mathcal{D}_s \cup (s_t, a_t, c_t, s_{t+1})$
 - 10: **else**
 - 11: $\mathcal{D} \leftarrow \mathcal{D} \cup (s_t, a_t, r_t, c_t, s_{t+1})$
 - 12: Train $\omega_{c_1}, \omega_{c_2}$ on $\mathcal{D} \cup \mathcal{D}_s$ (Eq. 2)
 - 13: Sample a batch of transitions $\mathcal{B} = \{(s, a, r, c, s')\} \in \mathcal{D}$
 - 14: Train $\omega_{r_1}, \omega_{r_2}$ using \mathcal{B} (Eq. 4)
 - 15: $\nu' \leftarrow \nu - \beta_\nu \nabla_\nu \mathcal{L}(\pi_\phi, \nu, \varepsilon, \alpha)$ using \mathcal{B} (Eq. 11)
 - 16: $\phi' \leftarrow \phi + \beta_\phi \nabla_\phi \mathcal{L}(\pi_\phi, \nu', \varepsilon, \alpha)$ using \mathcal{B} (Eq. 11)
 - 17: Resample $\mathcal{B}' = \{s \in \mathcal{D}\}$
 - 18: $\varepsilon' \leftarrow \varepsilon + \beta_\varepsilon \nabla_\varepsilon \mathcal{J}_\varepsilon(\pi_{\phi'})$ using \mathcal{B}' (Eq. 14)
 - 19: $\alpha' \leftarrow \alpha + \beta_\alpha \nabla_\alpha \mathcal{J}_\alpha(\pi_{\phi'}, \nu', \varepsilon')$ using \mathcal{D}_0 (Eq. 15)
 - 20: $\nu \leftarrow \nu', \phi \leftarrow \phi', \varepsilon \leftarrow \varepsilon', \alpha \leftarrow \alpha'$
 - 21: **if** $c_t == 1$ **then Break**
-

The objective function \mathcal{J}_ε is consistent with the objective function of the policy. This is evident by comparing Eq. 14 with the gradient w.r.t. the policy parameters ϕ in Eq. 11. The objective function for ε is designed to minimize the policy objective. This design stems from the idea that ε aims to capture the worst-case performance of the policy $\pi_{\phi'}$. Hence, by being optimized in this way, the safety region of the policy can be correctly adjusted to reflect that. It is important to note that we use ν'_{copy} to indicate that we merely use ν' value in the objective function and not include its gradient w.r.t. ε . We observed better performance by the gradient detachment in our early experiments which may be due to the injection of bias in ν' into its optimization process. Furthermore, to optimize the exploration value α , [9] used Q_{ω_r} as the objective function to change the value based on the performance of the policy. Therefore, in order to make the exploration rate of the Meta SAC-Lag safety compliant, we propose the objective function of α as:

$$\begin{aligned} \mathcal{J}_\alpha(\pi_{\phi'}, \nu', \varepsilon') &= \\ &= \max_{0 < \alpha \leq 1} \mathbb{E}_{s_0 \sim \rho_0} [Q_{\omega_r}(s_0, a) - \nu'(Q_{\omega_c}(s_0, a) - \varepsilon')] \end{aligned} \quad (15)$$

where $\pi_{\phi'}^{\text{det}}$ indicates the deterministic action value output by the policy. Basically, we use the expectation of the Lagrangian formulation evaluated in the initial states encoun-

tered by the agent. To gain a better understanding of the gradient relations, illustration of the optimization process of Meta SAC-Lag is depicted in Fig. 2.

D. Implementation Details

The learning process of Meta SAC-Lag is presented in Algorithm 1. The proposed algorithm utilizes three replay buffers for training. The main replay buffer \mathcal{D} stores all the transitions occurred while interacting with the environment, safety replay buffer \mathcal{D}_s stores all the transitions that have led to a constraint violation, and \mathcal{D}_0 builds an approximation of ρ_0 by generating samples from the distribution. We used a sampled batch $\mathcal{B} \subset \mathcal{D}$ to train the critic networks and the *inner* parameters ν and ϕ . Following that, as discussed in Section IV-A, we use resampled $\mathcal{B}' \subset \mathcal{D}$ and \mathcal{D}_0 to train the meta-parameters ε and α , respectively. The resampling process is analogous to the meta-testing process and is used to reduce bias in the training of the outer parameters [28], [29]. Moreover, following the original architecture [7], Meta SAC-Lag uses two critic and safety critic networks to prevent the overestimation of the value functions. To this end, the target values in Eq. 5 are calculated as $\min\{Q_{\bar{\omega}_{r_1}}, Q_{\bar{\omega}_{r_2}}\}$ and $\max\{Q_{\bar{\omega}_{c_1}}, Q_{\bar{\omega}_{c_2}}\}$, respectively. The $\bar{\omega}$ notation is used to indicate the target networks which are copies of the main networks updated with a time delay. Proposed in [30], the target networks aim to increase the stability of the training process and are calculated using the polyak averaging: $\bar{\omega} = \tau\omega + (1 - \tau)\bar{\omega}$. The hyperparameter $\tau \in (0, 1)$ typically has a value near zero.

Finally, it is also worth mentioning, in contrast to the original SAC, we use RMSProp [31] instead of Adam to calculate the higher-order gradients of the parameters ν , ϕ , and ε since backpropagating through RMSProp seems to be more numerically stable [9].

V. EXPERIMENTS

In this section, we evaluate the performance of Meta SAC-Lag. Specifically, our aim is to study two questions:

- How much does the added autonomy affect the performance of the algorithm compared to the baseline methods?
- How capable is Meta SAC-Lag to learn optimal performance in a real-world setup while avoiding actions that might catastrophically damage the system?

A. Test Benchmarks and Baselines

In order to study how the proposed algorithm will perform in safety-critical robotic scenarios, we use five simulated robotic environments with four different themes:

- **Locomotion:** In this theme, the purpose of control is to move the robotic system in the forward direction. The safety constraints are violated whenever the controller’s actions make the system exceed its limits, e.g., the velocity is higher than a certain threshold or the robot is falling to the ground. For that purpose, we use the Mujoco-based [32] Humanoid-Velocity environment from the Safety Gymnasium codebase [33]. It is important to note that the safety-related reward shaping of this

environment is removed to have a better understanding of the safety performance of the algorithms.

- **Obstacle Avoidance:** In many real-world robotic applications, there are mobile robots with manipulation capabilities. An important constraint of these systems is achieving their goal while avoiding certain regions in their surroundings. We adopt Isaac Gym-based Freight-FrankaCloseDrawer [34]. In this setup, the robot attempts to get near a drawer and close it while avoiding a red region. In addition, we use Car-Circle2 task [33] where the objective is to steer a car in a circular motion while avoiding collision with two walls.
- **Manipulation:** Another important area of safety-concerned robotic applications is manipulation. For that purpose, we use two embodied scenarios. For the **robotic manipulation** task we use Push Topple [35], [36] environment where the robotic arm must relocate a box without toppling it. Furthermore, in the **dexterous manipulation** scenario, we adopt the Egg Manipulate task where the agent must rotate an egg to a specific orientation without dropping it or exerting a force of more than 20 N. For both tasks, we use the Gymnasium Robotics codebase [37].

It is important to note that in the training process, we treat the constraints as hard constraints and terminate the episode whenever a violation has happened in the system. Furthermore, three baseline algorithms are chosen to compare and study the performance of Meta SAC-Lag:

- **SACv2-Lag:** The basic form of Meta SAC-Lag which uses Eq. 11 to optimize the policy and the Lagrangian multiplier with a fixed safety threshold.
- **Reward Constrained Policy Optimization (RCPO-SACv2):** Optimizes the policy using the Q -function formulated as $\mathbb{E}_\pi[\hat{Q}(s, a) = Q_r(s, a) - \nu Q_c(s, a)]$. The dual variable ν is also updated using Eq. 11.
- **RCPO-MetaSAC:** To show the effectiveness of our safe exploration technique, we use the $\hat{Q}(s, a)$ formulation in RCPO and optimize α using the approach proposed in [9].
- **Meta SAC-Lag \mathcal{J}_{nl} :** Inspired by [38], we experiment with a nonlinear objective function for ε specified as:

$$\mathcal{J}_\varepsilon^{nl}(\pi_{\phi'}) = \mathbb{E}_{\substack{s \sim \mathcal{D} \\ a \sim \pi_{\phi'}}} \left[\begin{cases} Q_{\omega_r}(s, a)Q_{\omega_c}(s, a) & \text{if } Q_{\omega_r}(s, a) < 0 \\ Q_{\omega_r}(s, a)(1 - Q_{\omega_c}(s, a)) & \text{otherwise} \end{cases} \right] \quad (16)$$

Essentially, $\mathcal{J}_\varepsilon^{nl}$ can have the advantage of no reliance on external parameter values, as opposed to Eq. 14 which uses ν' in the objective formulation.

To have a fair comparison, we tune the values of ε and ν for SACv2-Lag and RCPO-SACv2. Also, we use the values of RCPO-SACv2 for RCPO-MetaSAC. The values are outlined in Table I. For the value of α , SACv2 constrains the policy entropy as $\mathbb{E}_{s \sim \mathcal{D}}[-\log(\pi(s_t, a_t))] \geq \mathcal{H}$ and defines the α loss as $\min_{\alpha > 0} \mathcal{L}(\alpha) = \mathbb{E}_{s \sim \mathcal{D}}[\alpha(\log(\pi(s_t, a_t)) + \mathcal{H})]$. The authors propose the formula $\mathcal{H} = -\dim(\mathcal{A})$ as their target entropy. Furthermore, two important initial hyperparameter

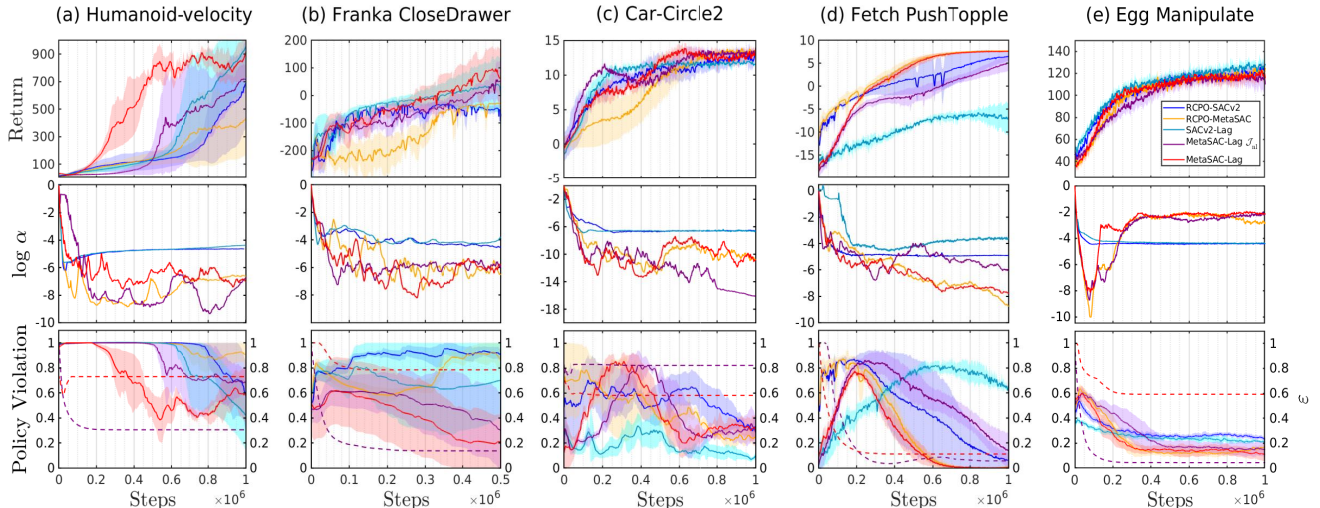


Fig. 3: Performance of Meta SAC-Lag compared with the baseline algorithms. **(Top row):** Reward performance during the learning process. (Higher values are better) **(Middle row):** The value of Exploration hyperparameter (α). **(Bottom row):** Episodic policy safety performance of the algorithms during the learning process. (Lower values are better). The dashed lines illustrate the constraint threshold value (ε).

TABLE I: Hyperparameter values (ε and ν) of the comparison methods

Environment / Parameter	ε	Meta SAC-Lag SACv2-Lag	RCPO-SACv2 RCPO-MetaSAC
Humanoid-Velocity	0.4	10	10
Franka DrawerClose	0.6	10	10
Car-Circle2	0.5	100	1
Fetch PushTopple	0.5	1000	10
Egg Manipulate	0.5	100	1

values of Meta SAC-Lag are automatically tuned; therefore, we set $\varepsilon = 1$ and $\alpha = 1$ as their initial values. Moreover, due to their similar training pipelines, we use the initial values of ν for SACv2-Lag in Table I for Meta SAC-Lag. We also set $\gamma_r = 0.99$ and $\gamma_c = 0.6$ for all the tasks. The results indicate the mean and variance of the performance of the algorithms across multiple independent runs.

B. Simulation Results

The simulation results are depicted in Fig. 3. To this end, we also report the return and the policy episodic violation rate in Table II. The violation rate is calculated as the average number of failures over a specific window of episodes. The results not only indicate that Meta SAC-Lag provides automated tuning of the safety-related hyperparameters but also, that the convergence process of the policy incurs lower constraint violations and yields higher or comparable returns. Furthermore, the update profile of α shows that as training goes on, in most cases, Meta SAC-Lag updates α to values lower than SACv2. This indicates that as the policy converges to a near-safe optimal solution, α is rapidly decreased to favor exploitation and prevent further constraint violations. Moreover, we can observe similar α profiles in Meta SAC-Lag and RCPO-SACv2 which can be attributed to α being optimized using similar objective functions. In addition, the optimization process of ε shows a generally

fast convergence. The fast convergence of ε provides the advantage of stable optimization as other values can be updated based on the optimally achieved value of ε . Finally, regarding the comparison between Eq. 14 and Eq. 16 we observe consistently better performance of Eq. 14 in both aspects of return and safety. In summary, the optimization outcomes of Meta SAC-Lag demonstrate that the algorithm excels across a range of embodied control tasks, proficiently learning optimal solutions, while demanding minimal hyperparameter tuning.

C. Real-World Deployment

Deployability can be regarded as one of the most important obstacles in using RL for learning to control real-world systems [39]. Choosing unsafe actions might lead the system to states that might damage it catastrophically, if chosen repeatedly. Therefore, using the conventional safe RL algorithms hinders their deployability since they require intensive hyperparameter tuning. In line with our purpose of assessing the deployability of a safe RL method, we propose a simple, yet important, safe RL testbench. This task, which we call *Pour Coffee*, is the task of moving a coffee-filled mug from a home position to a specific location and pouring the coffee into another cup. The task is executed using a Kinova Gen3 robot and its digital twin is created in the PyBullet simulation environment [40]. We define the state space $\mathcal{S} = \{X_{cup} \cup O_{cup} \cup \dot{X}_{cup} \cup X_{goal} \cup O_{goal}\}$ where $X = \{x, y, z\}$ and $O = \{\psi, \theta, \phi\}$ refer to the Cartesian position and the Euler angles in the Tait-Bryan ZYX intrinsic convention, respectively. Furthermore, the action of the agent maps to the velocity of the end-effector: $\mathcal{A} = \{\dot{x}_{cup}, \dot{y}_{cup}, \dot{z}_{cup}, \dot{\phi}_{cup}\}$. Moreover, we hierarchically define the reward function for reaching and pouring the coffee based on the Euclidean distance between the cup and

TABLE II: Relative Performance of the Algorithms During the Learning Process
(The best performance is shown in bold)

Task / Method	SACv2-Lag		RCPO-SACv2		RCPO-MetaSAC		MetaSAC-Lag		MetaSAC-Lag \mathcal{J}_{nl}	
	Jr	Jc	Jr	Jc	Jr	Jc	Jr	Jc	Jr	Jc
Humanoid-Velocity	956.11	0.42	690.66	0.57	435.93	0.91	812.90	0.49	712.72	0.67
Franka DrawerClose	2.19	0.69	-67.20	0.89	-29.85	0.85	63.13	0.20	54.02	0.30
Car-Circle2	11.65	0.10	12.67	0.32	12.36	0.24	13.79	0.31	13.95	0.29
Fetch-Topple	-6.6	0.66	6.55	0.06	7.6	0.007	7.49	0.004	4.98	0.17
Egg Manipulate	128.66	0.21	121.53	0.23	124.99	0.17	115.44	0.11	109.95	0.14

TABLE III: *Pour Coffee* Reward-Constraint Settings

Experiment Setting	Reward			Violation	
	Distance (r_1)	Acceleration (r_2)	Penalty (r_3)	Collision	Spillage
Simulation #1	✓	✗	✗	✓	✓
Simulation #2	✓	✓	✗	✓	✗
Simulation #3	✓	✓	✓	✓	✓
Simulation #4	✓	✓	✓	✓	✓
Real	✓	✗	✗	✓	✓

the goal $d = \|X_{cup} - X_{goal}\|_2$:

$$r(s, a, s') = \begin{cases} r_1 \cdot d + r_2 \cdot \|\ddot{X}_{cup}\| + r_3 \cdot \mathbb{1}[\text{spillage}] & \text{if } d > d_{\text{thresh}} \\ -|\phi_{cup} - \phi_{goal}| + 10 & \text{otherwise} \end{cases} \quad (17)$$

where $r_1 = -2$, $r_2 = -0.05$, $r_3 = -1$ and $d_{\text{thresh}} = 5$ cm. Furthermore, the system violates the safety constraints whenever self-collision or collision with the environment objects occurs. Additionally, we can define another constraint as spilling the coffee. As will be shown, this constraint forces the policy to be less jerky and aims to minimize the acceleration. The advantage of this approach, in contrast to similar environments [41], is the fact that it will eliminate the need to engineer the reward function to minimize the jerk and acceleration of the robot.

We conduct experiments with Meta SAC-Lag in four reward and constraint settings. The experiments aim to study whether formulating the problem sub-objectives can be more practical by defining them as constraints rather than shaping the reward explicitly. In the presented task, coffee spillage provides an implicit sub-objective that can be explicitly modeled as the sub-objective of minimizing the jerk and acceleration of the end-effector during the execution of the task. As shown in Table III, three experiments (Simulation #2, #3, #4) utilize different reward shaping schemes along with various constraint definitions. Moreover, we trained Meta SAC-Lag without engineered reward shaping (Simulation #1) both in the simulation environment and the real-world Kinova Gen3 setup. In order to make comparisons and evaluate the Sim2Real capability, the simulation-trained models were deployed on the robot using the checkpoints saved during the learning process.

The evaluation results are depicted in Fig. 4. The results illustrate that, as a result of providing a denser reward signal, explicit reward shaping can have positive effects in the increase of the success rate. However, using the spillage constraint helps the algorithm be even more effort-compliant resulting in lower jerk and comparable acceleration results. In other words, while being successful in executing the task is the most important metric, in a real-world scenario, sacrificing the performance to lower the effort of the system

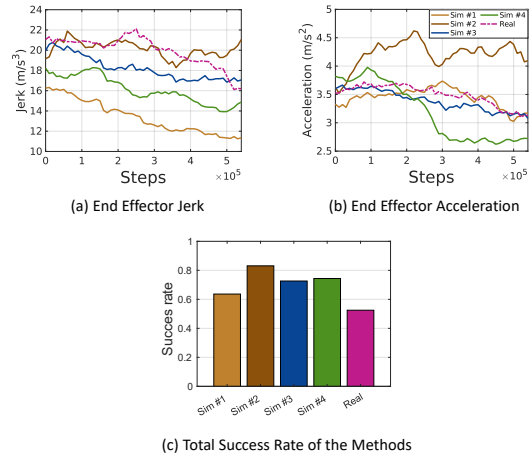


Fig. 4: Deployment results of Meta SAC-Lag on the real-world setup. (a) and (b) represent the jerk and acceleration of the end effector during the training process. (c) shows the final success rate of the algorithms.

and satisfy other safety concerns can be reasonable. In addition, regarding the comparison between Sim2Real and Real deployment of the proposed algorithm, we can observe that while both setups have similar behaviors, the real-world deployment is slightly hindered by the system’s physical limitations, such as sensor noise, control saturation, system fatigue, etc. Despite all that, the algorithm trained on the real-world setup without engineered reward function achieves results comparable to the models trained in the simulation.

VI. CONCLUSIONS

The paper focused on the problem of automatic hyperparameter tuning in Lagrangian safe RL methods. A novel model-free architecture called Meta SAC-Lag was proposed which addressed two inherent problems: safe exploration and constraint bound tuning. To this end, through the use of metagradient optimization, the algorithm is capable of adjusting the safety-related hyperparameters with minimal initial tuning. Furthermore, we studied the performance of our algorithm in five simulated embodied applications with the themes of locomotion, obstacle avoidance, robotic manipulation, and dexterous manipulation. We observed that the synergy created between the parameters and the hyperparameters results in comparable or better performance of the policy in terms of reward or safety. Additionally, we conducted an experiment in a real-world setup involving a practical coffee-pouring robotic environment without any explicit safety-related reward shaping. We deployed the algorithm on the Kinova Gen3 robot and showed that the

proposed algorithm can be helpful for real-world safety-sensitive applications by reducing the reliance on heuristic implementation of safety. We also observed that formulation of safety solely as the violation rather than engineering the reward function results in applying lower levels of effort at the cost of a diminished success performance. This trade-off can be especially favorable in real-world setups where safety violations are costly. Specifically, the proposed algorithm will learn the optimal policy in the real-world setup while adhering to the collision constraints and minimizing the effort imposed on the robot.

REFERENCES

- [1] A. Perera and P. Kamalaruban, "Applications of reinforcement learning in energy systems," *Renewable and Sustainable Energy Reviews*, vol. 137, p. 110618, 2021. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1364032120309023>
- [2] K. Shao, Z. Tang, Y. Zhu, N. Li, and D. Zhao, "A survey of deep reinforcement learning in video games," *arXiv preprint arXiv:1912.10944*, 2019.
- [3] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al., "Learning dexterous in-hand manipulation," *The International Journal of Robotics Research*, vol. 39, no. 1, pp. 3–20, 2020.
- [4] G. Dulac-Arnold, D. Mankowitz, and T. Hester, "Challenges of real-world reinforcement learning," *arXiv preprint arXiv:1904.12901*, 2019.
- [5] S. Gu, L. Yang, Y. Du, G. Chen, F. Walter, J. Wang, Y. Yang, and A. Knoll, "A review of safe reinforcement learning: Methods, theory and applications," *arXiv preprint arXiv:2205.10330*, 2022.
- [6] L. Brunke, M. Greeff, A. W. Hall, Z. Yuan, S. Zhou, J. Panerati, and A. P. Schoellig, "Safe learning in robotics: From learning-based control to safe reinforcement learning," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 5, pp. 411–444, 2022.
- [7] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," *arXiv preprint arXiv:1801.01290*, 2018.
- [8] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905*, 2019.
- [9] Y. Wang and T. Ni, "Meta-sac: Auto-tune the entropy temperature of soft actor-critic via metagradient," *arXiv preprint arXiv:2007.01932*, 2020.
- [10] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [11] E. Altman, *Constrained Markov decision processes*. Routledge, 2021.
- [12] D. Bertsekas, *Dynamic programming and optimal control: Volume I*. Athena scientific, 2012, vol. 4.
- [13] Y. Shen, M. J. Tobia, T. Sommer, and K. Obermayer, "Risk-sensitive reinforcement learning," *Neural Computation*, vol. 26, no. 7, pp. 1298–1328, 2014.
- [14] A. Stooke, J. Achiam, and P. Abbeel, "Responsive safety in reinforcement learning by pid lagrangian methods," in *International Conference on Machine Learning*. PMLR, 2020, pp. 9133–9143.
- [15] C. Tessler, D. J. Mankowitz, and S. Mannor, "Reward constrained policy optimization," *arXiv preprint arXiv:1805.11074*, 2018.
- [16] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *International conference on machine learning*. PMLR, 2017, pp. 1126–1135.
- [17] Z. Xu, H. P. van Hasselt, and D. Silver, "Meta-gradient reinforcement learning," *Advances in neural information processing systems*, vol. 31, 2018.
- [18] T. Haarnoja, S. Ha, A. Zhou, J. Tan, G. Tucker, and S. Levine, "Learning to walk via deep reinforcement learning," *arXiv preprint arXiv:1812.11103*, 2018.
- [19] V. Veeriah, T. Zahavy, M. Hessel, Z. Xu, J. Oh, I. Kemaev, H. P. van Hasselt, D. Silver, and S. Singh, "Discovery of options via meta-learned subgoals," *Advances in Neural Information Processing Systems*, vol. 34, pp. 29 861–29 873, 2021.
- [20] T. Zahavy, Z. Xu, V. Veeriah, M. Hessel, J. Oh, H. P. van Hasselt, D. Silver, and S. Singh, "A self-tuning actor-critic algorithm," *Advances in neural information processing systems*, vol. 33, pp. 20 913–20 924, 2020.
- [21] D. A. Calian, D. J. Mankowitz, T. Zahavy, Z. Xu, J. Oh, N. Levine, and T. Mann, "Balancing constraints and rewards with meta-gradient d4pg," *arXiv preprint arXiv:2010.06324*, 2020.
- [22] 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.
- [23] R. Koppejan and S. Whiteson, "Neuroevolutionary reinforcement learning for generalized control of simulated helicopters," *Evolutionary intelligence*, vol. 4, pp. 219–241, 2011.
- [24] G. Thomas, Y. Luo, and T. Ma, "Safe reinforcement learning by imagining the near future," *Advances in Neural Information Processing Systems*, vol. 34, pp. 13 859–13 869, 2021.
- [25] T. M. Moldovan and P. Abbeel, "Safe exploration in markov decision processes," *arXiv preprint arXiv:1205.4810*, 2012.
- [26] F. Berkenkamp, M. Turchetta, A. Schoellig, and A. Krause, "Safe model-based reinforcement learning with stability guarantees," *Advances in neural information processing systems*, vol. 30, 2017.
- [27] Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh, "A lyapunov-based approach to safe reinforcement learning," *Advances in neural information processing systems*, vol. 31, 2018.
- [28] A. Beirami, M. Razaviyayn, S. Shahrampour, and V. Tarokh, "On optimal generalizability in parametric learning," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [29] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil, "Bilevel programming for hyperparameter optimization and meta-learning," in *International conference on machine learning*. PMLR, 2018, pp. 1568–1577.
- [30] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [31] S. Ruder, "An overview of gradient descent optimization algorithms," *arXiv preprint arXiv:1609.04747*, 2016.
- [32] E. Todorov, T. Erez, and Y. Tassa, "Mujoco: A physics engine for model-based control," in *2012 IEEE/RSJ international conference on intelligent robots and systems*. IEEE, 2012, pp. 5026–5033.
- [33] J. Ji, B. Zhang, J. Zhou, X. Pan, W. Huang, R. Sun, Y. Geng, Y. Zhong, J. Dai, and Y. Yang, "Safety gymnasium: A unified safe reinforcement learning benchmark," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [34] J. Liang, V. Makoviychuk, A. Handa, N. Chentanez, M. Macklin, and D. Fox, "Gpu-accelerated robotic simulation for distributed reinforcement learning," in *Conference on Robot Learning*. PMLR, 2018, pp. 270–282.
- [35] H. Bharadhwaj, A. Kumar, N. Rhinehart, S. Levine, F. Shkurti, and A. Garg, "Conservative safety critics for exploration," *arXiv preprint arXiv:2010.14497*, 2020.
- [36] H.-L. Hsu, Q. Huang, and S. Ha, "Improving safety in deep reinforcement learning using unsupervised action planning," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 5567–5573.
- [37] R. de Lazcano, K. Andreas, J. J. Tai, S. R. Lee, and J. Terry, "Gymnasium robotics," 2023. [Online]. Available: <http://github.com/Farama-Foundation/Gymnasium-Robotics>
- [38] H. Honari, M. G. Tamizi, and H. Najjaran, "Safety optimized reinforcement learning via multi-objective policy optimization," *arXiv preprint arXiv:2402.15197*, 2024.
- [39] A. M. S. Enayati, R. Dershan, Z. Zhang, D. Richert, and H. Najjaran, "Facilitating sim-to-real by intrinsic stochasticity of real-time simulation in reinforcement learning for robot manipulation," *IEEE Transactions on Artificial Intelligence*, pp. 1–15, 2023.
- [40] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning," <http://pybullet.org>, 2016–2021.
- [41] Y. Zhu, J. Wong, A. Mandlekar, R. Martín-Martín, A. Joshi, S. Nasiriany, and Y. Zhu, "robosuite: A modular simulation framework and benchmark for robot learning," *arXiv preprint arXiv:2009.12293*, 2020.

A. Theoretical Analysis

In this section, the gradients of the each component of the algorithm is derived. While the automated differentiation tools for deep learning such as PyTorch and TensorFlow provide automated gradient calculation of this algorithm, the gradient analysis of Meta SAC-Lag may provide insightful information.

Step 1: In the beginning of optimization, the gradient of the Lagrangian multiplier ν is calculated using the inner loss (for ease of presentation, the expected values and the source of s are dropped in this analysis):

$$\begin{aligned}\nabla_{\nu}\mathcal{J}_{\nu} &= \nabla_{\nu}\mathcal{L}(\pi_{\phi}, \nu, \varepsilon, \alpha) \\ &= \nabla_{\nu} [Q_{\omega_r}(s, \pi_{\phi}(s)) - \nu(Q_{\omega_c}(s, \pi_{\phi}(s)) - \varepsilon) - \alpha \log \pi_{\phi}(a|s)] \\ &= -[Q_{\omega_c}(s, \pi_{\phi}(s)) - \varepsilon]\end{aligned}\quad (18)$$

Hence, ν is updated as:

$$\nu' \leftarrow \nu - \beta_{\nu}\nabla_{\nu}\mathcal{J}_{\nu} = \nu + \beta_{\nu} [Q_{\omega_c}(s, \pi_{\phi}(s)) - \varepsilon] \quad (19)$$

Step 2: Following the Lagrangian multiplier, the gradient of the actor parameters w.r.t. the Lagrangian loss is calculated as:

$$\begin{aligned}\nabla_{\phi}\mathcal{J}_{\phi} &= \mathcal{L}(\pi_{\phi}, \nu', \varepsilon, \alpha) \\ &= \nabla_{\phi} [Q_{\omega_r}(s, \pi_{\phi}(s)) - \nu'(Q_{\omega_c}(s, \pi_{\phi}(s)) - \varepsilon) - \alpha \log \pi_{\phi}(a|s)] \\ &= \nabla_{\phi} [Q_{\omega_r}(s, \pi_{\phi}(s)) - (\nu + \beta_{\nu}(Q_{\omega_c}(s, \pi_{\phi}(s)) - \varepsilon))(Q_{\omega_c}(s, \pi_{\phi}(s)) - \varepsilon) - \alpha \log \pi_{\phi}(a|s)] \\ &= \nabla_{\phi}Q_{\omega_r}(s, \pi_{\phi}(s)) - \nabla_{\phi}Q_{\omega_c}(s, \pi_{\phi}(s)) [2\beta_{\nu}Q_{\omega_c}(s, \pi_{\phi}(s)) - 2\beta_{\nu}\varepsilon + \nu] - \alpha\nabla_{\phi} \log \pi_{\phi}(a|s)\end{aligned}\quad (20)$$

The actor parameters are then updated as:

$$\begin{aligned}\phi' &\leftarrow \phi + \beta_{\phi}\nabla_{\phi}\mathcal{J}_{\phi} \\ &= \phi + \beta_{\phi} [\nabla_{\phi}Q_{\omega_r}(s, \pi_{\phi}(s)) - \nabla_{\phi}Q_{\omega_c}(s, \pi_{\phi}(s)) [2\beta_{\nu}Q_{\omega_c}(s, \pi_{\phi}(s)) - 2\beta_{\nu}\varepsilon + \nu] \\ &\quad - \alpha\nabla_{\phi} \log \pi_{\phi}(a|s)]\end{aligned}\quad (21)$$

Step 3: The first meta-parameter used in the training pipeline is the safety threshold ε . Assuming the meta-objective formulation of Equation 14, the gradient of ε can be derived as:

$$\begin{aligned}\nabla_{\varepsilon}\mathcal{J}_{\varepsilon} &= \nabla_{\varepsilon} [\nu'_{\text{copy}}Q_{\omega_c}(s, \pi_{\phi'}(s)) - Q_{\omega_r}(s, \pi_{\phi'}(s))] \\ &= \nabla_{\varepsilon}\phi'^{\mathbf{T}}\nabla_{\phi'}\mathcal{J}_{\varepsilon} + \nabla_{\varepsilon}\nu' \left[\nabla_{\nu'}\phi'^{\mathbf{T}}\nabla_{\phi'}\mathcal{J}_{\varepsilon} + \cancel{\nabla_{\nu'}\mathcal{J}_{\varepsilon}} \right] + \cancel{\nabla_{\varepsilon}\mathcal{J}_{\varepsilon}}\end{aligned}\quad (22)$$

It should be noted that the chain of gradients that are calculated must be causal. For that reason, gradients such as $\nabla_{\phi'}\nu'$ would be meaningless and have not been written. Moreover, each component of Equation 22 can be computed as:

$$\nabla_{\varepsilon}\phi' = 2\beta_{\phi}\beta_{\nu}\nabla_{\phi}Q_{\omega_c}(s, \pi_{\phi}(s)) \quad (23)$$

$$\nabla_{\phi'}\mathcal{J}_{\varepsilon} = \nu'\nabla_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}(s)) - \nabla_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}(s)) \quad (24)$$

$$\nabla_{\varepsilon}\nu' = -\beta_{\nu} \quad (25)$$

$$\nabla_{\nu'}\phi' = -\beta_{\phi}\nabla_{\phi}Q_{\omega_c}(s, \pi_{\phi}(s)) \quad (26)$$

Therefore, the final gradient is calculated as:

$$\nabla_{\varepsilon}\mathcal{J}_{\varepsilon} = [-3\beta_{\nu}\beta_{\phi}\nabla_{\phi}Q_{\omega_c}(s, \pi_{\phi}(s))]^{\mathbf{T}} [\nabla_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}(s)) - \nu'\nabla_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}(s))] \quad (27)$$

Hence, the value of ε is updated as:

$$\begin{aligned}\varepsilon' &\leftarrow \varepsilon + \beta_{\varepsilon}\nabla_{\varepsilon}\mathcal{J}_{\varepsilon} \\ &= \varepsilon + \beta_{\varepsilon} [-3\beta_{\nu}\beta_{\phi}\nabla_{\phi}Q_{\omega_c}(s, \pi_{\phi}(s))]^{\mathbf{T}} [\nabla_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}(s)) - \nu'\nabla_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}(s))]\end{aligned}\quad (28)$$

The meta-gradient of ε consists of two components. While the right hand side evaluates the performance of the new actor parameter ϕ' by calculating the tradeoff between safety and optimality, the left hand side is the main direction driver of the gradient. The left hand side of the meta-gradient equation determines the policy's degree of unsafety with respect to the safety critic. The gradient of the policy w.r.t. the safety critic would determine the effect of the changes in the policy parameters on the measure of safety. The gradient will always try to force the policy to become safer by moving in the

negative direction of the gradient. The idea of minimizing the policy objective discussed in Section IV-C is also clear from the observation of the gradients of ε and specifically in Equation 24 and the left hand side of Equation 27. Without the negative sign the algorithm would try to match ε with the unsafety performance of the policy rather than forcing it to become safer.

Step 4: The final step of the optimization involves calculating the meta-gradient of the temperature α which is calculated as:

$$\begin{aligned}\nabla_{\alpha}\mathcal{J}_{\alpha} &= \nabla_{\alpha}[Q_{\omega_r}(s, \pi_{\phi'}^{det}(s)) - \nu'(Q_{\omega_c}(s, \pi_{\phi'}^{det}(s)) - \varepsilon')] \\ &= \cancel{\nabla_{\alpha}\nu'}[\nabla_{\nu'}\mathcal{J}_{\alpha} + \nabla_{\nu'}\phi'^{\mathbf{T}}\nabla_{\phi'}\mathcal{J}_{\alpha} + \nabla_{\nu'}\varepsilon'\nabla_{\varepsilon'}\mathcal{J}_{\alpha}] \\ &\quad + \nabla_{\alpha}\phi'^{\mathbf{T}}[\nabla_{\phi'}\mathcal{J}_{\alpha} + \nabla_{\phi'}\varepsilon'\nabla_{\varepsilon'}\mathcal{J}_{\alpha}] + \cancel{\nabla_{\alpha}\varepsilon'}\nabla_{\varepsilon'}\mathcal{J}_{\alpha}\end{aligned}\tag{29}$$

The components of the gradient can then be calculated as:

$$\nabla_{\varepsilon'}\mathcal{J}_{\alpha} = \nu' \tag{30}$$

$$\nabla_{\alpha}\phi' = -\beta_{\phi}\nabla_{\phi}\log\pi_{\phi}(a|s) \tag{31}$$

$$\nabla_{\phi'}\mathcal{J}_{\alpha} = \nabla_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}^{det}(s)) - \nu'\nabla_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}^{det}(s)) \tag{32}$$

$$\nabla_{\phi'}\varepsilon' = \beta_{\varepsilon}[-3\beta_{\nu}\beta_{\phi}\nabla_{\phi}Q_{\omega_c}(s, \pi_{\phi}(s))]^{\mathbf{T}}[H_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}(s)) - \nu'H_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}(s))] \tag{33}$$

Therefore, the final gradient of α is obtained as:

$$\begin{aligned}\nabla_{\alpha}\mathcal{J}_{\alpha} &= -\beta_{\phi}\nabla_{\phi}\log\pi_{\phi}(a|s)^{\mathbf{T}}[\nabla_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}^{det}(s)) - \nu'\nabla_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}^{det}(s))] \\ &\quad + \nu'\beta_{\varepsilon}[-3\beta_{\nu}\beta_{\phi}\nabla_{\phi}Q_{\omega_c}(s, \pi_{\phi}(s))]^{\mathbf{T}}[H_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}(s)) - \nu'H_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}(s))] \\ &= -\beta_{\phi}\nabla_{\phi}\log\pi_{\phi}(a|s)^{\mathbf{T}}[\nabla_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}^{det}(s)) - \nu'\nabla_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}^{det}(s))]\end{aligned}\tag{34}$$

The second term of the gradient ($\nabla_{\phi'}\varepsilon'\nabla_{\varepsilon'}\mathcal{J}_{\alpha}$) was dropped because in the architecture of the neural networks used to implement Meta SAC-Lag, the ReLU activation function was used. An important feature of ReLU is the fact that it has zero second-order derivative almost everywhere. Hence, the Hessian matrix in Equation 33 would be equal to zero.

Finally, the update rule for ε follows:

$$\begin{aligned}\alpha' &\leftarrow \alpha + \beta_{\alpha}\nabla_{\alpha}\mathcal{J}_{\alpha} \\ &= \alpha - \beta_{\alpha}\beta_{\phi}\nabla_{\phi}\log\pi_{\phi}(a|s)^{\mathbf{T}}[\nabla_{\phi'}Q_{\omega_r}(s, \pi_{\phi'}^{det}(s)) - \nu'\nabla_{\phi'}Q_{\omega_c}(s, \pi_{\phi'}^{det}(s))]\end{aligned}\tag{35}$$

By observing the right hand side of $\nabla_{\alpha}\mathcal{J}_{\alpha}$ it can be noticed that term is equivalent to $\nabla_{\phi'}[Q_{\omega_r}(s, \pi_{\phi'}^{det}(s)) - \nu'Q_{\omega_c}(s, \pi_{\phi'}^{det}(s))]$ which is similar to the critic formulation of RCPO as $\hat{Q}^{\pi}(s, a) = Q^{\pi}(s, a) - \nu Q_c^{\pi}(s, a)$. Hence, in essence the meta-gradient calculation of α in Meta SAC-Lag and RCPO-MetaSAC are the same. This observation can be confirmed by noticing the similarity in the updating profile of α in the simulation results in Figure 3. The profile is also similar to that of Meta SAC-Lag \mathcal{J}_{nl} because due to the use of ReLU and the Hessian matrix becoming zero, the effect of $\mathcal{J}_{\varepsilon}$ on α has been eliminated and the updating procedure for both of them would be the same.