# MOSEAC: Streamlined Variable Time Step Reinforcement Learning

Dong Wang

dong-1.wang@polymtl.ca Department of Computer Engineering Polytechnique Montréal Giovanni Beltrame

giovanni.beltrame@polymtl.ca Department of Computer Engineering Polytechnique Montréal

# Abstract

Traditional reinforcement learning (RL) methods typically employ a fixed control loop, where each cycle corresponds to an action. This rigidity poses challenges in practical applications, as the optimal control frequency is task-dependent. A suboptimal choice can lead to high computational demands and reduced exploration efficiency. Variable Time Step Reinforcement Learning (VTS-RL) addresses these issues by using adaptive frequencies for the control loop, executing actions only when necessary. This approach, rooted in reactive programming principles, reduces computational load and extends the action space by including action durations. However, VTS-RL's implementation is often complicated by the need to tune multiple hyperparameters that govern exploration in the multi-objective action-duration space (i.e., balancing task performance and number of time steps to achieve a goal). To overcome these challenges, we introduce the Multi-Objective Soft Elastic Actor-Critic (MOSEAC) method. This method features an adaptive reward scheme that adjusts hyperparameters based on observed trends in task rewards during training. This scheme reduces the complexity of hyperparameter tuning, requiring a single hyperparameter to guide exploration, thereby simplifying the learning process and lowering deployment costs. We validate the MOSEAC method through simulations in a Newtonian kinematics environment, demonstrating high task and training performance with fewer time steps, ultimately lowering energy consumption. This validation shows that MOSEAC streamlines RL algorithm deployment by automatically tuning the agent control loop frequency using a single parameter. Its principles can be applied to enhance any RL algorithm, making it a versatile solution for various applications.

#### 1 Introduction

Model-free deep reinforcement learning (RL) algorithms have achieved significant success in various domains, including gaming (Liu et al., 2021; Ibarz et al., 2021) and robotic control (Akalin & Loutfi, 2021; Singh et al., 2022). Traditional RL methods typically rely on a fixed control loop where the agent takes actions at predetermined intervals, e.g. every 0.1 seconds (10 Hz). However, this approach introduces significant limitations: fixed-rate control can lead to stability issues and high computational demands, especially in dynamic environments where the optimal action frequency may vary over time.

To address these limitations, Variable Time Step Reinforcement Learning (VTS-RL) was recently introduced, based on reactive programming principles (Majumdar et al., 2021; Bregu et al., 2016). The key idea behind VTS-RL is to *execute control actions only when necessary*, thereby reducing computational load, as well as expanding the action space including variable action durations. For instance, in an autonomous driving scenario, VTS-RL enables the agent to adapt its control strategy

based on the situation, using a low frequency in stable conditions and a high frequency in more complex situations (Wang & Beltrame, 2024a).

Two notable VTS-RL algorithms are Soft Elastic Actor-Critic (SEAC) (Wang & Beltrame, 2024a) and Continuous-Time Continuous-Options (CTCO) (Karimi et al., 2023), which simultaneously learn actions and their durations.

CTCO employs continuous-time decision-making and flexible option durations, improving exploration and robustness in continuous control tasks. However, its effectiveness can be hampered by the need to tune several hyperparameters, and the fact that tasks can take a long time to complete as CTCO does not optimize for task duration.

In contrast, SEAC adds reward terms for task energy (i.e. the numbers of actions performed) and task time (i.e. the time needed to complete a task), showing effectiveness in time-restricted environments like a racing video game (Wang & Beltrame, 2024b). However, SEAC requires careful tuning of its hyperparameters (balancing task, energy, and time costs) to avoid performance degradation.

SEAC and CTCO's brittleness to hyperparameter settings poses a challenge for users aiming to fully leverage their potential. To mitigate this issue, we propose the Multi-Objective Soft Elastic Actor-Critic (MOSEAC) algorithm based on SEAC (Wang & Beltrame, 2024a). MOSEAC adds action durations to the action space and adjusts hyperparameters based on observed trends in task rewards during training, leaving a single hyperparameter to be set to guide the exploration. We evaluate MOSEAC in a Newtonian kinematics simulation environment, comparing it CTCO, SEAC, and fixed-frequency SAC, showing that MOSEAC is stable, accelerates training, and has high final task performance. Additionally, our hyperparameter setting approach can be broadly applied to any continuous action reinforcement learning algorithm, such as Twin Delayed Deep Deterministic (TD3) (Fujimoto et al., 2018) or Proximal Policy Optimization (PPO) (Schulman et al., 2017). This adaptability facilitates the transition from fixed-time step to variable-time step reinforcement learning, significantly expanding the scope of its practical application.

# 2 Related Work

The importance of action duration in reinforcement learning algorithms has been severely underestimated for a long time. In fact, if reinforcement learning algorithms are to be applied to the real world, it is an essential factor, impacting an agent's exploration capabilities (Amin et al., 2020; Park et al., 2021). High frequencies may reduce exploration efficiency in some cases, but are needed in others such as environments with delays (Bouteiller et al., 2021). Wang & Beltrame (2024b); Karimi et al. (2023) show how various frequencies impact learning, with high frequencies potentially impeding convergence. Thus, a dynamically optimal control frequency, which adjusts in real time based on reactive principles, could enhance performance and adaptability.

Sharma et al. (2017) introduced a concept similar to variable control frequencies, developing a framework for repetitive action reinforcement learning. This model allows an agent to perform identical actions over successive states, combining them into a larger action. This method has intrigued researchers, especially in gaming (Metelli et al., 2020; Lee et al., 2020), but fails to account for physical properties or reduce computational demands, making its practical application challenging.

Additionally, Chen et al. (2021) attempted to modify the traditional "control rate" by integrating actions such as "sleep" to reduce activity periods ostensibly. However, this approach still mandates fixed-frequency checks of system status, which does not effectively diminish the computational load as anticipated. These instances underscore the ongoing challenges and the nascent stage of effectively integrating variable control frequencies and repetitive behaviors into real-world applications, highlighting a critical gap between theoretical innovation and practical efficacy.

# 3 Algorithm Framework

To reduce hyperparameter dimensions, we combine SEAC's (Wang & Beltrame, 2024a) hyperparameters balancing task, energy, and time rewards using a simple multiplication method. We then apply an adaptive adjustment method on the remaining hyperparameters. Like SEAC, our reward equation includes components for:

- Quality of task execution (the standard RL reward),
- Time required to complete a task (important for varying action durations), and
- Energy (the number of time steps, a.k.a. the number of actions taken, which we aim to minimize).

#### Definition 1

$$R = \alpha_m R_t R_\tau - \alpha_t$$

where  $R_t$  is the task reward;  $R_{\tau}$  is a time dependent term.

 $\alpha_m$  is a weighting factor used to modulate the magnitude of the reward. Its primary function is to prevent the reward from being too small, which could lead to task failure, or too large, which could cause reward explosion, ensuring stable learning.

 $\alpha_{\varepsilon}$  is a penalty parameter applied at each time step to impose a cost on the agent's actions. This parameter gives a fixed cost to the execution of an action, thereby discouraging unnecessary ones. In practice  $\alpha_{\varepsilon}$  promotes the completion of a task using fewer time steps (remember that time steps have variable duration), i.e. it reduces the energy used by the control loop of the agent.

We determine the optimal policy  $\pi^*$ , which maximizes the reward R.

The reward as designed, minimizes both energy cost (number of steps) and the total time to complete the task through  $R_{\tau}$ . By scaling the task-specific reward based on action duration with  $\alpha_m$ , agents are motivated to complete tasks using fewer actions:

**Definition 2** The remap relationship between action duration and reward

$$R_{\tau} = t_{min}/t, \quad R_{\tau} \in [t_{min}/t_{max}, 1]$$

where t is the duration of the current action,  $t_{min}$  is the minimum duration of an action (strictly greater than 0), and  $t_{max}$  is the maximum duration of an action.

To automatically set  $\alpha_m$  and  $\alpha_{\varepsilon}$  to optimal values, we adjust them dynamically during training. Based on Wang & Beltrame (2024b)'s experience, we increase  $\alpha_m$  and decrease  $\alpha_{\varepsilon}$  to mitigate convergence issues, specifically the problem of sparse rewards caused by a suboptimal set of a large  $\alpha_{\varepsilon}$  and a small  $\alpha_m$ . This adjustment ensures that rewards are appropriately balanced, facilitating learning and convergence.

These parameters are adjusted based on observed trends in task rewards: if the average reward is declining (see Definition 4), we increase  $\alpha_m$  and decrease  $\alpha_{\varepsilon}$ , linking them with a sigmoid function, ensuring a balanced reward structure that promotes convergence and efficient learning. To guarantee stability, we ensure the change is monotonic, forcing a uniform sweep of the parameter space.

**Definition 3** The relationship between  $\alpha_{\varepsilon}$  and  $\alpha_m$  is

$$\alpha_{\varepsilon} = 0.2 \cdot \left(1 - \frac{1}{1 + e^{-\alpha_m + 1}}\right)$$

Based on Wang & Beltrame (2024a)'s experience, we establish a mapping relationship between the two parameters: when the initial value of  $\alpha_m$  is 1.0, the initial value of  $\alpha_{\varepsilon}$  is 0.1. As  $\alpha_m$  increases,  $\alpha_{\varepsilon}$  decreases, but never falls below 0.

To determine the trend on the average reward, we perform a linear regression across the current training episode and compute the slope of the resulting line: if it is negative, the reward is declining.

**Definition 4** The slope of the average reward  $(k_R)$  is:

$$k_R(R_a) = \frac{n \sum_{i=1}^n (i \cdot R_{ai}) - (\sum_{i=1}^n i) (\sum_{i=1}^n R_{ai})}{n \sum_{i=1}^n i^2 - (\sum_{i=1}^n i)^2}$$

where n is the total number of data points collected across the update interval  $(k_{update} \text{ in Algorithm } 1)$ . The update interval is a hyperparameter that determines the frequency of updates for these neural networks used in the actor and critic policies, occurring after every n episode (Sutton & Barto, 2018).  $R_a$  represents the list of average rewards  $(R_{a1}, R_{a2}, ..., R_{an})$  during training. Here, an average reward  $R_{ai}$  is calculated across one episode.

After observing a downward trend in the average reward during the training process, we introduced the hyperparameter  $\psi$  to dynamically adjust  $\alpha_m$  as defined in Definition 5.

**Definition 5** We adaptively adjust the reward every  $k_{update}$  episodes. When  $k_R < 0$ :

 $\alpha_m = \alpha_m + \psi$ 

where,  $\psi$  is the only additional hyperparameter required by our MOSEAC to tune the reward equation during training. Additionally,  $\alpha_{\varepsilon}$  is adjusted as described in Definition 3.

Algorithm 1 shows the pseudocode of MOSEAC. Convergence analysis is provided in Appendix A.

Algorithm 1: Multi-Objective Soft Elastic Actor and Critic

**Require:** a policy  $\pi$  with a set of parameters  $\theta$ ,  $\theta'$ , critic parameters  $\phi$ ,  $\phi'$ , variable time step environment model  $\Omega$ , learning-rate  $\lambda_p$ ,  $\lambda_q$ , reward buffer  $\beta_r$ , replay buffer  $\beta$ . 1: Initialization  $i = 0, t_i = 0, \beta_r = 0$ , observe  $S_0$ while  $t_i \leq t_{max}$  do 2: for  $i \leq k_{length} \vee Not Done \mathbf{do}$ 3:  $A_i, D_i = \pi_{\theta}(S_i)$  $\rightarrow$  simple action and its duration 4:  $S_{i+1}, R_i = \Omega(A_i, D_i)$  $\rightarrow$  compute reward with Definition 1 and 2 5:  $i \leftarrow i + 1$ 6: end 7:  $\beta_r \leftarrow 1/i \times \sum_{i=1}^{i} R_i$  $\rightarrow$  collect the average reward for one episode 8:  $\beta \leftarrow S_{0 \sim i}, A_{0 \sim i}, D_{0 \sim i}, R_{0 \sim i}, S_{1 \sim i+1}$ 9: i = 010:  $t_i \leftarrow t_i + 1$ 11: 12:if  $t_i \ge k_{init}$  &  $t_i \mid k_{update}$  then Sample S, A, D, R, S' from( $\beta$ ) 13: $\phi \leftarrow \phi - \lambda_q \nabla_\delta \mathcal{L}_Q(\phi, S, A, D, R, S')$  $\rightarrow$  critic update 14:  $\theta \leftarrow \theta - \lambda_p \nabla_{\theta} \mathcal{L}_{\pi}(\theta, S, A, D, \phi)$  $\rightarrow$  actor update 15:if  $k_R(\beta_r)$  then 16: $\alpha_m \leftarrow \alpha_m + \psi$  $\rightarrow$  see Definition 4 for  $k_R$ 17: $\alpha_{\varepsilon} \leftarrow F_{update}(\alpha_m)$ 18:  $\rightarrow$  update  $\alpha_m, \alpha_{\varepsilon}$  followed Definition 3 19: end  $\beta_r = 0$  $\rightarrow$  Re-record average reward values under new hyperparameters 20: end 21: Perform soft-update of  $\phi'$  and  $\theta'$ 22: 23: end

Here,  $t_{max}$  represents maximum training steps (Sutton & Barto, 2018);  $k_{length}$  is maximum exploration steps per episode (Sutton & Barto, 2018);  $k_{init}$  is steps in the initial random exploration

phase (Sutton & Barto, 2018). The reward  $R_i$  is calculated as  $R(S_i, A_i, D_i)$ , where  $D_i$  lies within  $[T_{min}, T_{max}]$ , representing action duration.

Our algorithm optimizes multiple objectives using a streamlined, weighted strategy. Unlike Hierarchical Reinforcement Learning (HRL) (Dietterich, 2000; Li et al., 2019), which seeks Pareto optimality (Kacem et al., 2002; Monfared et al., 2021) with layered reward policies, our method simplifies the approach. We emphasize user-friendliness and computational efficiency, making our strategy easily adaptable to various algorithms.

Apart from a series of hyperparameters inherent to RL that need adjustment, such as learning rate,  $t_{max}$ ,  $k_{update}$ , etc.,  $\psi$  is the only hyperparameter that requires tuning in MOSEAC.

However, one limitation of our approach is that high  $\psi$  values can lead to issues such as reward explosion. While a small  $\psi$  value generally avoids significant problems and guides  $\alpha_m$  to its optimal value, it requires extended training periods. Determining the appropriate  $\psi$  value remains a critical optimization point in our method. We recommend using the pre-set  $\psi$  value provided in our implementation, thus reducing the need for additional parameter adjustments. If users encounter problems such as reward explosion, it is advisable to reduce the  $\psi$  value appropriately. Our ongoing research aims to mitigate further the risks associated with  $\psi$ , enhancing the algorithm's reliability and robustness.

### 4 Empirical Analysis

We conduct six experiments involving MOSEAC, SEAC, CTCO, and SAC algorithms (at 20 Hz and 60 Hz) within a Newtonian kinematics environment. A detailed description of the environment appears in Appendix B. These tests are performed on a computer equipped with an Intel Core i5-13600K and a Nvidia RTX 4070 GPU. The operating system is Ubuntu 22.04 LTS. For the hyperparameter settings of MOSEAC, please refer to Appendix C<sup>1</sup>

Figure 1 illustrates the result of training process. In Figure 1a, the average returns align with our expectations, showing that CTCO is sensitive to the choice of action duration range, which affects its discount factor  $\gamma$ . During the training process  $\gamma$  tended to be very small, compromising CTCO's long-term planning. Overall, CTCO performs worse than the two fixed-control frequency SAC baselines in this environment. In contrast, MOSEAC trains slightly faster than SEAC, as shown in Figure 1b.





(a) Average returns of 5 reinforcement learning algorithms in 3M steps during the training.

(b) Average Energy cost of 5 reinforcement learning algorithms in 3M steps during the training.

Figure 1: Result graphs of five reinforcement learning algorithms.

<sup>&</sup>lt;sup>1</sup>Our code is public available.

We compared MOSEAC's task performance with the best-performing SEAC after training. Due to the poor performance of CTCO and SAC, we excluded them from the analysis. Figure 2 presents performance metrics related to energy consumption and task duration. Specifically, Figure 2a shows the average energy consumption (measured in steps) for 300 tasks, and Figure 2b compares average task durations. Data show that MOSEAC's average energy and time costs are lower than SEAC's.

Wilcoxon signed-rank test (chosen due to the lack of normality in the data distribution) indicates that MOSEAC's energy cost is lower than SEAC's (z = -1.823, p = 0.020), where MOSEAC has a mean energy cost of 3.120 (SD = 0.424), compared to SEAC's mean of 3.193 (SD = 0.451). Detailed statistical results are provided in Appendix D. Similarly, MOSEAC's time cost is significantly lower than SEAC's (z = -2.669, p = 0.004), with MOSEAC having a mean time cost of 0.905 (SD = 0.125), compared to SEAC's mean of 0.935 (SD = 0.127).

The improved performance of MOSEAC over SEAC can be attributed to its reward function. While SEAC's reward function is linear, combining task reward, energy penalty, and time penalty independently, MOSEAC introduces a multiplicative relationship between task reward and time-related reward. This non-linear interaction enhances the reward signal, particularly when both task performance and time efficiency are high, and naturally balances these factors. By keeping the energy penalty separate, MOSEAC maintains flexibility in tuning without complicating the relationship between time and task rewards. This design allows MOSEAC to more effectively guide the agent's decisions, resulting in better energy efficiency and task completion speed in practical applications.



(a) The energy cost (counted by numbers of steps) figure for 300 different tasks.

(b) The time cost figure for 300 different tasks.

Figure 2: Comparison figure of energy consumption and time performance of 300 different tasks. We use the same random seed to ensure that both have the same tasks.

# 5 Conclusion and Future Work

In this paper, we present MOSEAC, a VTS-RL algorithm with adaptive hyperparameters that respond to observed reward trends during training. This approach significantly reduces hyperparameter sensitivity and increases robustness, improving data efficiency. Unlike other VTS-RL algorithms, our method does not requires a single hyperparameter, thereby reducing learning and tuning costs when transitioning from fixed to variable time steps. This establishes a strong foundation for realworld applications.

Our next objective is to implement our method in real-world robotics applications, such as smart cars and robotic arms.

### References

- Neziha Akalin and Amy Loutfi. Reinforcement learning approaches in social robotics. Sensors, 21 (4):1292, 2021.
- Susan Amin, Maziar Gomrokchi, Hossein Aboutalebi, Harsh Satija, and Doina Precup. Locally persistent exploration in continuous control tasks with sparse rewards. *arXiv preprint arXiv:2012.13658*, 2020.
- Yann Bouteiller, Simon Ramstedt, Giovanni Beltrame, Christopher Pal, and Jonathan Binas. Reinforcement learning with random delays. In *International conference on learning representations*, 2021.
- Endri Bregu, Nicola Casamassima, Daniel Cantoni, Luca Mottola, and Kamin Whitehouse. Reactive control of autonomous drones. In *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services*, pp. 207–219, 2016.
- Yuxin Chen, Hejun Wu, Yongheng Liang, and Guoming Lai. Varlenmarl: A framework of variablelength time-step multi-agent reinforcement learning for cooperative charging in sensor networks. In 2021 18th Annual IEEE International Conference on Sensing, Communication, and Networking (SECON), pp. 1–9. IEEE, 2021.
- Thomas G Dietterich. Hierarchical reinforcement learning with the maxq value function decomposition. Journal of artificial intelligence research, 13:227–303, 2000.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actorcritic methods. In International conference on machine learning, pp. 1587–1596. PMLR, 2018.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. arXiv preprint arXiv:1812.05905, 2018a.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. arXiv preprint arXiv:1812.05905, 2018b.
- Julian Ibarz, Jie Tan, Chelsea Finn, Mrinal Kalakrishnan, Peter Pastor, and Sergey Levine. How to train your robot with deep reinforcement learning: lessons we have learned. *The International Journal of Robotics Research*, 40(4-5):698–721, 2021.
- Imed Kacem, Slim Hammadi, and Pierre Borne. Pareto-optimality approach for flexible job-shop scheduling problems: hybridization of evolutionary algorithms and fuzzy logic. *Mathematics and* computers in simulation, 60(3-5):245–276, 2002.
- Amirmohammad Karimi, Jun Jin, Jun Luo, A Rupam Mahmood, Martin Jagersand, and Samuele Tosatto. Dynamic decision frequency with continuous options. In 2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 7545–7552. IEEE, 2023.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- Vijay Konda and John Tsitsiklis. Actor-critic algorithms. Advances in neural information processing systems, 12, 1999.
- Jongmin Lee, Byung-Jun Lee, and Kee-Eung Kim. Reinforcement learning for control with multiple frequencies. Advances in Neural Information Processing Systems, 33:3254–3264, 2020.
- Siyuan Li, Rui Wang, Minxue Tang, and Chongjie Zhang. Hierarchical reinforcement learning with advantage-based auxiliary rewards. Advances in Neural Information Processing Systems, 32, 2019.

- Rongrong Liu, Florent Nageotte, Philippe Zanne, Michel de Mathelin, and Birgitta Dresp-Langley. Deep reinforcement learning for the control of robotic manipulation: a focussed mini-review. *Robotics*, 10(1):22, 2021.
- Rupak Majumdar, Aman Mathur, Marcus Pirron, Laura Stegner, and Damien Zufferey. Paracosm: A test framework for autonomous driving simulations. In *International Conference on Fundamental Approaches to Software Engineering*, pp. 172–195. Springer International Publishing Cham, 2021.
- Alberto Maria Metelli, Flavio Mazzolini, Lorenzo Bisi, Luca Sabbioni, and Marcello Restelli. Control frequency adaptation via action persistence in batch reinforcement learning. In *International Conference on Machine Learning*, pp. 6862–6873. PMLR, 2020.
- Mohammadali Saniee Monfared, Sayyed Ehsan Monabbati, and Atefeh Rajabi Kafshgar. Paretooptimal equilibrium points in non-cooperative multi-objective optimization problems. *Expert Systems with Applications*, 178:114995, 2021.
- Seohong Park, Jaekyeom Kim, and Gunhee Kim. Time discretization-invariant safe action repetition for policy gradient methods. Advances in Neural Information Processing Systems, 34:267–279, 2021.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
- Sahil Sharma, Aravind Srinivas, and Balaraman Ravindran. Learning to repeat: Fine grained action repetition for deep reinforcement learning. arXiv preprint arXiv:1702.06054, 2017.
- Bharat Singh, Rajesh Kumar, and Vinay Pratap Singh. Reinforcement learning in robotic applications: a comprehensive survey. Artificial Intelligence Review, 55(2):945–990, 2022.
- Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT press, 2018.
- Dong Wang and Giovanni Beltrame. Deployable reinforcement learning with variable control rate. arXiv preprint arXiv:2401.09286, 2024a.
- Dong Wang and Giovanni Beltrame. Reinforcement learning with elastic time steps. arXiv preprint arXiv:2402.14961, 2024b.

# A Convergence Analysis of MOSEAC

This appendix aims to analyze the convergence of the Multi-Objective Soft Elastic Actor-Critic (MOSEAC) algorithm. Its action space includes a time dimension D and the reward function is  $R = \alpha_m R_t R_\tau - \alpha_\varepsilon$ , with  $\alpha_m$  and  $\alpha_\varepsilon$  dynamically adjusted during training, details can be found in section 3.

#### A.1 SAC Algorithm Overview

In the standard SAC algorithm (Haarnoja et al., 2018a), the policy  $\pi_{\theta}(a|s)$  selects action a, and updates the policy parameters  $\theta$  and value function parameters  $\phi$ . The objective function is:

$$J(\pi_{\theta}) = \mathbb{E}_{(s,a) \sim \pi_{\theta}} \left[ Q^{\pi}(s,a) + \alpha \mathcal{H}(\pi_{\theta}(\cdot|s)) \right]$$

where  $J(\pi_{\theta})$  is the objective function,  $Q^{\pi}(s, a)$  is the state-action value function,  $\alpha$  is a temperature parameter controlling the entropy term, and  $\mathcal{H}(\pi_{\theta}(\cdot|s))$  is the entropy of the policy.

#### A.2 Incorporating Time Dimension and Our Reward Function

When the action space is extended to include D and the reward function is modified to  $R = \alpha_m R_t R_\tau - \alpha_\varepsilon$ , where  $\alpha_m \ge 0$ ,  $0 < R_\tau \le 1$ , and  $\alpha_\varepsilon$  is a small positive constant, the new objective function becomes:

$$J(\pi_{\theta}) = \mathbb{E}_{(s,a,D)\sim\pi_{\theta}} \left[ Q^{\pi}(s,a,D) + \alpha \mathcal{H}(\pi_{\theta}(\cdot|s)) \right]$$

where  $Q^{\pi}(s, a, D)$  is the extended state-action value function with time dimension D, and  $R_t$  and  $R_{\tau}$  are components of the reward function, see Definition 1.

#### A.3 Policy Gradient

With our reward function, the policy gradient is:

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(a, D|s) \left( Q^{\pi}(s, a, D) \cdot (\alpha_{m} \cdot R_{\tau}) - \alpha_{\varepsilon} \right) \right]$$

where  $\nabla_{\theta} J(\pi_{\theta})$  is the gradient of the objective function with respect to the policy parameters  $\theta$ .

#### A.4 Value Function Update

The value function update, incorporating the time dimension D and our reward function, is:

$$L(\phi) = \mathbb{E}_{(s,a,D,r,s')} \left[ \left( Q_{\phi}(s,a,D) - \left( r + \gamma \mathbb{E}_{(a',D')\sim\pi_{\theta}} \left[ V_{\phi}(s') - \alpha \log \pi_{\theta}(a',D'|s') \right] \right) \right)^2 \right]$$

where  $L(\phi)$  is the loss function for the value function update, r is the reward,  $\gamma$  is the discount factor, and  $V_{\bar{\phi}}(s')$  is the target value function.

#### A.5 Policy Parameter Update

The new policy parameter  $\theta$  update rule is:

$$\theta_{k+1} = \theta_k + \beta_k \mathbb{E}_{s \sim D, (a, D) \sim \pi_{\theta}} \left[ \nabla_{\theta} \log \pi_{\theta}(a, D|s) \left( Q_{\phi}(s, a, D) \cdot (\alpha_m \cdot R_{\tau}) - \alpha_{\varepsilon} - V_{\bar{\phi}}(s) + \alpha \log \pi_{\theta}(a, D|s) \right) \right]$$

where  $\beta_k$  is the learning rate at step k.

#### A.6 Dynamic Adjustment and Convergence Analysis

To analyze the impact of dynamically adjusting  $\alpha_m$  and  $\alpha_{\varepsilon}$ , we assume:

#### 1. Dynamic Adjustment Rules:

-  $\alpha_m$  increases monotonically by a small increment  $\psi$  if the reward trend decreases over consecutive episodes (See Definition 4). -  $\alpha_{\varepsilon}$  decreases with the Definition 3.

#### 2. Learning Rate Conditions:

The learning rates  $\alpha_k$  and  $\beta_k$  must satisfy the following conditions (Konda & Tsitsiklis, 1999):

$$\sum_{k=0}^{\infty} \alpha_k = \infty, \quad \sum_{k=0}^{\infty} \alpha_k^2 < \infty$$
$$\sum_{k=0}^{\infty} \beta_k = \infty, \quad \sum_{k=0}^{\infty} \beta_k^2 < \infty$$

#### A.7 Unbiased Estimation

Assuming the critic estimates are unbiased:

$$\mathbb{E}[Q_{\phi}(s, a, D) \cdot (\alpha_m \cdot R_{\tau}) - \alpha_{\varepsilon}] = Q^{\pi}(s, a, D) \cdot (\alpha_m \cdot R_{\tau}) - \alpha_{\varepsilon}$$

Since  $R_{\tau}$  is a positive number within [0, 1], its effect on  $Q^{\pi}(s, a, D)$  is linear and does not affect the consistency of the policy gradient.

#### A.8 Impact Evaluation

#### 1. Positive Scaling:

As  $R_{\tau}$  is always a positive number less than or equal to 1, and  $\alpha_m \geq 0$ , it only scales the reward, not altering its sign. This scaling does not change the direction of the policy gradient but affects its magnitude.

#### 2. Small Offset:

 $\alpha_{\varepsilon}$  is a small constant used to accelerate training. This small offset does not affect the direction of the policy gradient but introduces a minor shift in the value function, which does not alter the overall policy update direction.

#### A.9 Conclusion

Under the aforementioned conditions, the modified SAC algorithm with the time dimension D and the new reward function  $R = \alpha_m R_t R_\tau - \alpha_\varepsilon$ , which we called MOSEAC, will converge to a local optimum, i.e. (Sutton & Barto, 2018),

$$\lim_{k\to\infty}\nabla_\theta J(\pi_\theta) = 0$$

This implies that over time, the policy parameters will stabilize at a local optimum, maximizing the policy performance.

# **B** Details of Newtonian kinematics environment

To validate our algorithm, we set up an autonomous driving environment based on Newtonian dynamics, including a random birth point, a random endpoint, and a random obstacle. Check Definition 7 for its state value information, Definition 8 for its action value information, and Definition 6 for the physical formula it follows.

As Definition 1, the precise reward configuration for our environment are outlined in Table 1.

**Definition 6** The Newton Kinematics formulas are:

$$\begin{split} D_{aim} &= 1/2 \cdot (V_{aim} + V_{current}) \cdot T; \\ V_{aim} &= V_{current} + AT; \\ F_{aim} &= mA; \\ F_{true} &= F_{aim} - f_{friction}; \end{split}$$

 $f_{friction} = \mu mg, \text{ if } F_{aim} > f_{friction} \wedge V_{agent} \neq 0$  $f_{friction} = F_{aim}, \text{ if } F_{aim} \leq f_{friction} \wedge V_{agent} = 0$ 

where  $D_{aim}$  is the distance generated by the policy that the agent needs to move.  $V_{agent}$  is the speed of the agent. T is the time to complete the movement generated by the policy,  $F_{aim}$  is the traction force necessary for the agent to change to the aimed speed. m is the mass of the agent,  $\mu$  is the friction coefficient, and g is the acceleration of gravity.

Table 1: Reward Settings for The Newtonian Kinematics Environment

		Reward Settings
Name	Value	Annotation
	500.0	Reach the goal
r	-500.0	Crash on an obstacle
	$D_{origin} - 1.0 \times D_{goal}$	$D_{origin}$ : distance from start to goal; $D_{goal}$ : distance to goal
$\alpha_m$	1.0	

The state dimensions are shown in Definition 7:

**Definition 7** The positions of the endpoint and obstacle are randomized in each episode.

 $S_t = (Pos, Obs, Goal, Speed, Time, Force)$ 

Where:

Pos = Position Data of The Agent in X and Y Direction,

Obs = Position Data of The Obstacle in X and Y Direction,Goal = Position Data of The Goal in X and Y Direction,

Speed = Spped Data of The Agent for Current Step,

Speed = Spped Data of The Agent for Current Step,

Time = The Duration data of The Agent to Execute The Last Time Step,

Force = Force Data of The Agent for Last Time Step in X and Y Direction.

The Spatial Information of our environment are shown in Table 2:

The action dimensions are shown in Definition 8:

Environment details				
Name	Value	Annotation		
Action dimension	3	,		
Range of speed	[-2,2]	m/s		
Action Space	[-100.0, 100.0]	Newton		
Range of time	[0.01, 1.0]	second		
State dimension	11	Task gain factor		
World size	(2.0, 2.0)	in meters		
Obstacle shape	Round	Radius: 5cm		
Agent weight	20	in $Kg$		
Gravity factor	9.80665	in $m/s^2$		
Static friction coeddicient	0.6			

Table 2.	Details of	The Simpl	o Nowtonian	Kinomatics (	Cymnasium	Environment
Table 2.	Detans of	. incompi	c rew toman	<b>I</b> III autos	Gymnasium.	Linvironniciti

# **Definition 8** Action

Each action should have a corresponding execution duration.

$$a_t = (D_t, A_{fx}, A_{fy})$$

Where:

 $D_t = The \ duration \ of \ the \ Agent \ to \ implement \ current \ action,$  $A_{fx} = The \ Force \ of \ the \ Agent \ in \ X \ Coordinate,$  $A_{fy} = The \ Force \ of \ the \ Agent \ in \ Y \ Coordinate.$ 

# C Hyperparameters of MOSEAC

Hyperparameter sheet of MOSEAC					
Name	Value	Annotation			
Total steps	3e6				
$\gamma$	0.99	Discount factor			
Net shape	(256, 256)				
batch_size	256				
a_lr	3e - 5	Learning rate of Actor Network			
c_lr	3e - 5	Learning rate of Critic Network			
max_steps	500	Maximum steps for one episode			
$\alpha$	0.12				
$\eta$	-3	Refer to SAC (Haarnoja et al., 2018b)			
min_time	0.01	Minimum control duration, in seconds			
max_time	1.0	Maximum control duration, in seconds			
$lpha_m$	1.0	Init value of $\alpha_m$ (Definition 1)			
$\psi$	1e - 4	Monotonically increasing H-parameter (Algorithm 1)			
Optimizer	Adam	Refer to Adam (Kingma & Ba, 2014)			
environment steps	1				
Replaybuffer size	1e6				
Number of samples before training start	$5 \cdot max\_steps$				
Number of critics	2				

# D Statistic Results

Measure 1		Measure 2	W	Z	р
MOSEAC Energy Cost MOSEAC Time Cost	-	SEAC Energy Cost SEAC Time Cost	2080.500 18561.000	-1.823 -2.669	$0.020 \\ 0.004$

Table 3: Paired Samples T-Test

rabic 1. rest of rollinging (Shapiro Vink)
--

			W	р
MOSEAC Energy Cost	-	SEAC Energy Cost	$0.775 \\ 0.996$	< 0.001
MOSEAC Time Cost	-	SEAC Time Cost		0.714

	Ν	Mean	SD	SE	Coefficient of variation
MOSEAC Energy Cost	300	3.120	0.424	0.024	0.136
SEAC Energy Cost	300	3.193	0.451	0.026	0.141
MOSEAC Time Cost	300	0.905	0.125	0.007	0.139
SEAC Time Cost	300	0.935	0.127	0.007	0.136