

# Enhancing Information Freshness: An AoI Optimized Markov Decision Process Dedicated in The Underwater Task

Yimian Ding<sup>\*,+</sup>, Jingzehua Xu<sup>\*,+</sup>, Yiyuan Yang<sup>†</sup>, Guanwen Xie<sup>\*</sup>, Shuai Zhang<sup>‡</sup>

<sup>\*</sup>MicroMasters Program in Statistics and Data Science, Massachusetts Institute of Technology, USA

<sup>†</sup>Department of Computer Science, University of Oxford, United Kingdom

<sup>‡</sup>Department of Data Science, New Jersey Institute of Technology, USA

Email: sz457@njit.edu

**Abstract**—Ocean exploration utilizing autonomous underwater vehicles (AUVs) via reinforcement learning (RL) has emerged as a significant research focus. However, underwater tasks have mostly failed due to the observation delay caused by acoustic communication in the Internet of underwater things. In this study, we present an AoI optimized Markov decision process (AoI-MDP) to improve the performance of underwater tasks. Specifically, AoI-MDP models observation delay as signal delay through statistical signal processing, and includes this delay as a new component in the state space. Additionally, we introduce wait time in the action space, and integrate AoI with reward functions to achieve joint optimization of information freshness and decision-making for AUVs leveraging RL for training. Finally, we apply this approach to the multi-AUV data collection task scenario as an example. Simulation results highlight the feasibility of AoI-MDP, which effectively minimizes AoI while showcasing superior performance in the task. To accelerate relevant research in this field, we have made the simulation codes available as open-source<sup>1</sup>.

**Index Terms**—Age of Information, Markov Decision Process, Statistical Signal Processing, Reinforcement Learning, Autonomous Underwater Vehicles

## I. INTRODUCTION

Harsh ocean environment put forward higher difficulty on ocean exploration [1]. As a novel approach, utilizing autonomous underwater vehicles (AUVs) via reinforcement learning (RL) has merged as a significant research focus [2]. Relying on Internet of underwater things (IoUT) [3], AUVs can communicate with each other and work in collaboration to accomplish human-insurmountable tasks [4]. However, underwater tasks have mostly failed due to the observation delay caused by acoustic communication, leading to the non-causality of control policies [5]. Although this issue can be alleviated by introducing states that incorporate past information and account for the future effects of control laws [5], it becomes increasingly challenging as the number of AUVs grows, leading to more complexity in both communication and decision-making processes [6].

As a significant indicator evaluation the freshness of information, age of information (AoI) is proposed to measure the

time elapsed at the receiver since the last information was generated until the most recent information is received [7]. And it has been verified to solve the severe delay caused by constantly sampling and transmitting observation information [8]. Central to this consensus is that minimizing AoI can enhance the freshness of information, thereby facilitating efficiency of subsequent decision-making process in the presence of observation delay [8]. Currently, numerous studies have focused on optimizing AoI to aid decision-making in the context of land-based or underwater tasks. For example, Messaoudi *et al.* optimized vehicle trajectories relying on minimizing average AoI while reducing energy consumption [9]. Similarly, Lyu *et al.* leveraged AoI to assess transmission delay impacts on state estimation, improving performance under energy constraints [10]. These studies primarily aim to reduce AoI by improving motion strategies of agents, without considering the impact of information update strategies on AoI. They assume that agents instantaneously perform the current action upon receiving previous information. However, this zero-wait strategy has been shown to be suboptimal in scenarios with high variability in delay times [11]. Conversely, it has been demonstrated that introducing waiting time before updating can achieve lower average AoI. This highlights the necessity of integrating optimized information update strategies into underwater tasks.

Furthermore, most studies currently leverage the standard Markov decision process (MDP) without observation delay to model the underwater tasks, which assumes the AUV can instantaneously receive current state information without delay, so that it can make corresponding actions [12]. This idealization, however, may not hold in many practical scenarios, since signal propagation delays and high update frequencies causing channel congestion reduce the freshness of received information, hindering the AUV's decision-making efficiency. Therefore, extending the standard MDP framework to incorporate observation delays and AoI is necessary [13].

Based on the above analysis, we attempted to propose an AoI optimized MDP (AoI-MDP) dedicated in the underwater task to improve the performance of the tasks with observation delay. The contributions of this paper include the following:

- To the best of our knowledge, we are the first to formulate

<sup>+</sup> These authors contributed equally to this work.

<sup>1</sup> The source code associated with this article is available at the following GitHub repository: <https://github.com/Xiboxtg/AoI-MDP>.

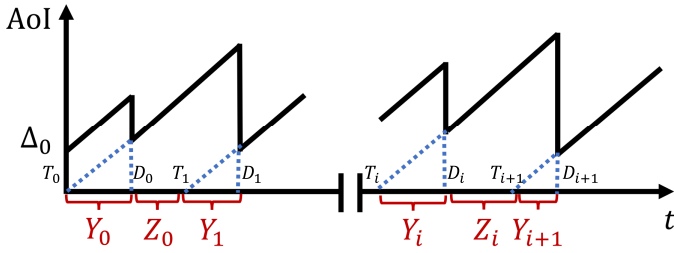


Fig. 1: Illustration of the AoI model, which is defined using a sawtooth piecewise function, where  $Y_i$  and  $Z_i$  denote the observation delay and wait time at time  $i$ , respectively.

the underwater task as an MDP that incorporates observation delay and AoI. Based on AoI-MDP, we utilize RL for AUV training to realize joint optimization of both information updating and decision-making strategies.

- Instead of simply modeling observation delay as a random distribution or stationary stochastic variable, we utilize statistical signal processing to realize the high-precision modeling via AUV equipped sonar, which potentially yielding more realistic results.
- Through comprehensive evaluations and ablation experiments in the underwater data collection task, our AoI-MDP showcases superior feasibility and excellent performance in balancing multi-objective optimization. And to accelerate relevant research in this field, the code for simulation will be released as open-source in the future.

## II. METHODOLOGY

In this section, we present the proposed AoI-MDP, which consists of three main components: an observation delay-aware state space, an action space that incorporates wait time, and reward functions based on AoI. To achieve high-precision modeling, AoI-MDP utilizes statistical signal processing (SSP) to represent observation delay as underwater acoustic signal delay, thereby aiming to minimize the gap between simulation and real-world underwater tasks.

### A. AoI Optimized Markov Decision Process

As illustrated in Fig. 1, consider the scenario where the  $i$ -th underwater acoustic signal is transmitted from the AUV at time  $T_i$ , and the corresponding observed information is received at time  $D_i$ , AoI is defined using a sawtooth piecewise function

$$\Delta(t) = t - T_i, D_i \leq t < D_{i+1}, \forall i \in \mathbf{N}. \quad (1)$$

Hence, we denote the MDP that integrates observation delay and characterizes the freshness of information through the AoI as the AoI-MDP, which can be defined by a quintuple  $\Omega$  for further RL training [14]

$$\Omega = \{\mathcal{S}, \mathcal{A}, \mathcal{R}, \Pr(s_{i+1}|s_i, a_i), \gamma\}, \quad (2)$$

where  $\mathcal{S}, \mathcal{A}, \mathcal{R}$  represent state space, action space and reward functions, respectively. The term  $\Pr(s_{i+1}|s_i, a_i) \in [0, 1]$  indicates state transition probability distribution, while  $\gamma \in [0, 1]$  represents a discount factor.

In AoI-MDP, instead of simply incorporating AoI as a component of the reward functions to guide objective optimization through RL training, we also leverage AoI as crucial side information to facilitate decision making. Specifically, we reformulate the standard MDP's state space, action space, and reward functions. The detailed designs for each of these elements are as follows:

**State Space  $\mathcal{S}$ :** the state space of the AoI-MDP consists of two parts: AUV's observed information  $s'_i$ , and observation delay  $Y_i$  at time  $i$ , represented by  $s_i = (s'_i, Y_i) \in \mathcal{S}' \times \mathcal{Y}$ . We introduce the observation delay  $Y_i$  as a new element so that the AUV can be aware of the underwater acoustic signal delay when the sonar emits an underwater acoustic signal to detect the surrounding environment. Additionally, we achieve high-precision modeling of both  $s'_i$  and  $Y_i$  through SSP, whose details are presented in Section II-B.

**Action Space  $\mathcal{A}$ :** the action space of the AoI-MDP consists of the tuple  $a_i = (a'_i, Z_i) \in \mathcal{A}' \times \mathcal{Z}$ , where  $a'_i$  denotes the actions taken by the AUV, while  $Z_i$  indicates the wait time between observing the environmental information and decision-making at time  $i$ . Through jointly optimizing the wait time  $Z_i$  and action  $a_i$ , we aim to minimize the AoI, enabling the AUV's decision-making policy to converge to an optimal level.

**Reward Function  $\mathcal{R}$ :** the reward function  $r'_i$  in standard MDP comprises elements with different roles, such as penalizing failures, promoting efficiency, and encouraging cooperation, etc. Here, we introduce the time-averaged AoI as a new component of the reward function. Thus, the updated reward function can be represented by the tuple  $r_i = (r'_i, -\bar{\Delta})$ . And the time-averaged AoI can be computed as follows:

$$\bar{\Delta} = \frac{\sum_{i=1}^{\mathcal{N}} ((2Y_{i-1} + Y_i + Z_{i-1}) \times (Y_i + Z_{i-1})) + S_0}{2 \times (\sum_{i=1}^{\mathcal{N}} Z_{i-1} + \sum_{i=1}^{\mathcal{N}} Y_i + Y_0)}, \quad (3)$$

where  $\mathcal{N}$  is the length of information signal,  $S_0 = 0.5 \times (2\Delta_0 + Y_0) \times Y_0$ . Therefore, the time averaged AoI can be minimized through RL training.

According to the above analysis, the total reward function is set below:

$$R_i = \sum_{k=1}^{\infty} \lambda^{(k)} r_i^{(k)}, \quad (4)$$

where  $\lambda^{(k)}$  represents the weighting coefficient of the  $k$ -th reward function.

Based on proposed AoI-MDP, we further integrate it with RL training for the joint optimization of information freshness and decision-making for AUVs. The pseudocode for the AoI-MDP based RL training is showcased in Algorithm 1.

### B. Observation Delay and Information Modeling

Different from previous work, our study enhances the state space of AoI-MDP by considering the observed information using estimated information perceived by AUV-equipped sensors. And we consider observation delay as underwater acoustic signal delay, rather than merely treating it as a random distribution [11] or stationary stochastic variable [15], [16].

---

**Algorithm 1: AoI-MDP Based RL Training**


---

```

1 Initialize the replay buffer  $\mathcal{D}$ , critic network, and actor
  network parameters of each AUV.
2 for each epoch  $k$  do
3   Reset the training environment and parameters.
4   for each time step  $i$  do
5     for each AUV  $j$  do
6       Obtain current state  $s'_i$  and observation
        delay  $Y_i$ .
7       Sample action  $a'_i$  and waiting time  $Z_i$ 
        according to the actor network.
8       Wait  $Z_i$  and execute  $a'_i$  while receiving
        reward  $R_i$ .
9       while In delay period do
10        Extract  $N$  tuples of data
          $(s_n, a_n, R_n, s_{n+1})_{n=1, \dots, N}$  from  $\mathcal{D}$ .
11        Update the Critic Network.
12        Update the Actor Network.
13      end
14      Store  $(s_i, a_i, R_i, s_{i+1})$  in  $\mathcal{D}$ .
15      while In waiting period do
16        Repeat the process in delay period.
17      end
18    end
19  end
20 end

```

---

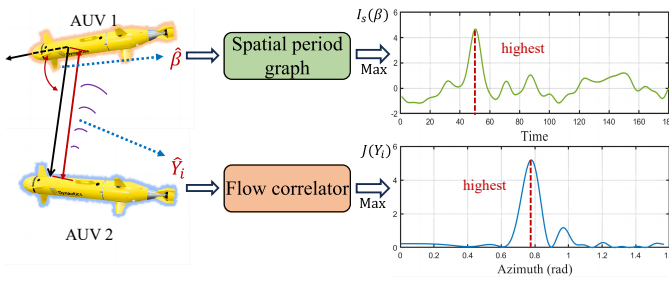


Fig. 2: Illustration of the azimuth and time delay estimation.

This approach aims to provide high-precision modeling to improve the performance in the underwater environment. And the schematic diagram is shown in Fig. 2.

To be specific, our study assumes the AUV leverages a sonar system to estimate the distance from itself to environmental objects. This was achieved by transmitting acoustic signals through sonar, measuring the time delay taken for these signals to propagate to the target, reflect, and return to the hydrophone, thus allowing for distance estimation. The acoustic signal propagation can be represented as

$$\mathcal{X}[n] = \mathcal{S}[n - Y_i] + \mathcal{W}[n], n = 0, 1, \dots, N - 1, \quad (5)$$

where  $\mathcal{S}[n]$  represents the known signal, while  $Y_i$  denotes the time delay to be estimated, and  $\mathcal{W}[n]$  is the Gaussian white noise with variance  $\sigma^2$ .

We further employ the flow correlator as an estimator to determine the time delay. Specifically, this estimator carries out the following computations on each received signal:

$$J[Y_i] = \sum_{n=Y_i}^{Y_i+M-1} \mathcal{X}[n] \mathcal{S}[n - Y_i], 0 \leq Y_i \leq N - M, \quad (6a)$$

$$\hat{Y}_i = \operatorname{argmax} [J[Y_i]], \quad (6b)$$

where  $M$  is the sampling length of  $\mathcal{S}[n]$ . By calculating the value of  $\hat{Y}_i$  that maximizes the value of Eq. (6a), the estimated time delay value can be obtained through Eq. (6b).

On the other hand, the AUV in our study utilizes a long linear array sensor to estimate the azimuth  $\beta$  between its orientation and environmental objects. The signal propagation can be expressed as follows:

$$x[n] = A \cos \left[ 2\pi \left( F_0 \frac{d}{c} \cos \beta \right) n + \phi \right] + \mathcal{W}[n], \quad (7)$$

$$n = 0, 1, \dots, M - 1,$$

where  $F_0$  denotes the frequency of transmitted signal, while  $d$  represents the interval between sensors. Besides,  $c$  indicates the speed of underwater acoustic signal propagation, while  $A$  and  $\phi$  are unknown signal amplitude and phase, respectively.

The estimator in SSP is further leveraged to estimate the azimuth  $\beta$ . By maximizing the spatial period graph, the estimate of  $\beta$  ( $0 < \beta < \pi/2$ ) can be calculated

$$I_s(\beta) = \frac{1}{M} \left( \left| \sum_{n=0}^{M-1} x[n] \exp[-j2\pi(F_0 \frac{d}{c} \cos \beta)n] \right| \right)^2, \quad (8a)$$

$$\hat{\beta} = \operatorname{argmax} [I_s(\beta)]. \quad (8b)$$

By calculating the value of  $\beta$  that maximizes the value of Eq. (8a), the estimated time delay value can be obtained through Eq. (8b).

Finally, the AUV can achieve target positioning using the observed  $\hat{Y}_i$  and  $\hat{\beta}$ . These observations are then utilized as data for the observed information in the state space of the AoI-MDP, which potentially yields more realistic results to improve the underwater performance, while reducing the gap between simulation and reality in the underwater tasks.

### III. EXPERIMENTS

In this section, we validate the proposed AoI-MDP through extensive simulation experiments. Further, we present the experimental results with further analysis and discussion.

#### A. Task Description and Settings

Since open-source underwater tasks are scarce, we consider the scenario of a multi-AUV data collection task as a classic example to evaluate the feasibility and effectiveness of the AoI-MDP. This task utilizes RL algorithms to train AUVs to collect data of sensor nodes in the Internet of underwater things, encompassing multiple objectives, such as maximizing sum data rate and collision avoidance, while minimizing energy consumption, etc. For the remaining details and parameters of the task, please refer to the previous work [4].

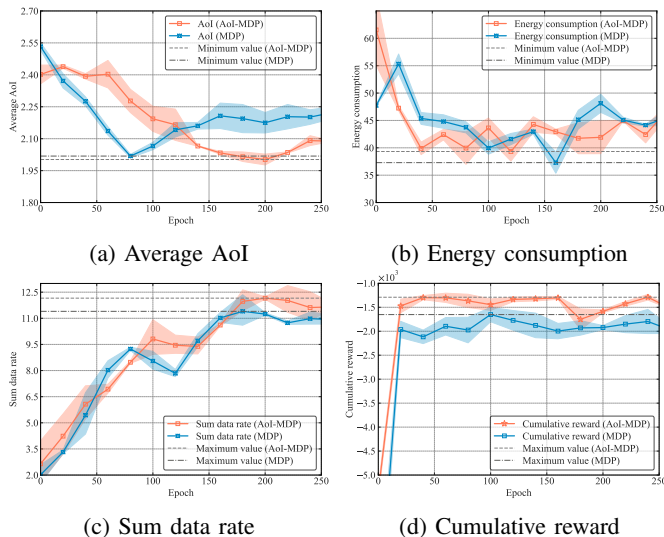


Fig. 3: Comparison of experimental results of RL training based on AoI-MDP and standard MDP.

TABLE I: Comparison of different delay models.

	AoI	Sum data rate	Energy consumption
SSP	1.97±0.26	11.99±0.73	33.83±2.59
Poisson	3.42±0.18	5.95±2.42	34.27±7.98
Exponential	2.67±0.26	7.65±1.99	43.11±4.16
Geometric	2.38±0.28	12.34±0.79	58.15±9.49

### B. Experiment Results and Analysis

We first compared the experimental results of RL training based on AoI-MDP and standard MDP under identical conditions, respectively. Results in Fig. 3 show that AoI-MDP results in lower time-averaged AoI, reduced energy consumption, higher sum data rate, and greater cumulative rewards. This demonstrates that AoI-MDP improves the training effectiveness and performance of the RL algorithm.

Then we evaluated the generalization performance of AoI-MDP using commonly employed delay models in the communication field, including exponential, poisson and geometric distributions. The experimental results, compared with the SSP model, are shown in Table 1. The AoI-MDP based RL training demonstrates superior performance across various distributions, indicating strong generalization capabilities. Additionally, SSP for time delay modeling achieved near-optimal results in AoI optimization, sum data rate optimization, and energy consumption optimization, underscoring its effectiveness in the underwater data collection task.

We further turned our attention to comparing the generalization of AoI-MDP in various RL algorithms. We conducted experiments utilizing AoI-MDP on soft actor-critic (SAC) and conservative Q-learning (CQL), within the contexts of online and offline RL, respectively. As shown in Fig. 4, both online and offline RL algorithms can successfully adapt to AoI-MDP, while ultimately achieving favorable training outcomes.

Finally we guided the multi-AUV in the underwater data

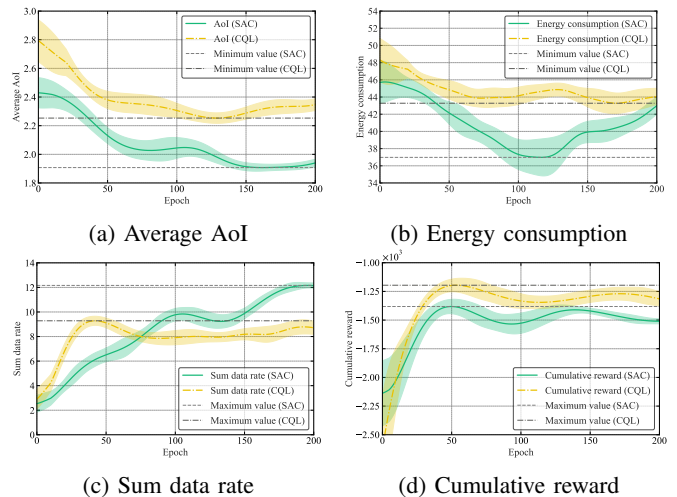
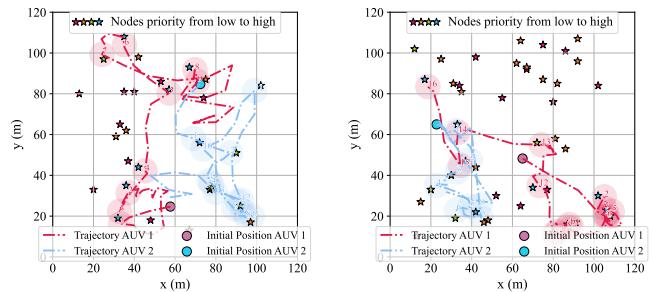


Fig. 4: Comparison of experimental results using online and offline RL algorithms based on AoI-MDP.



(a) AUV trajectories (AoI-MDP) (b) AUV trajectories (MDP)

Fig. 5: The AUV trajectories using the expert policy trained via SAC algorithm based on AoI-MDP and standard MDP.

collection task using the expert policy trained via SAC algorithm based on AoI-MDP and standard MDP respectively. As illustrated in Fig. 5, the trajectory coverage trained under AoI-MDP is more extensive, leading to more effective completion of the data collection task. Conversely, under standard MDP, the trajectories of AUVs appears more erratic, with lower node coverage, thereby showcasing suboptimal performance.

## IV. CONCLUSION

In this study, we propose AoI-MDP to improve the performance of underwater tasks. AoI-MDP models observation delay as signal delay through SSP, and includes this delay as a new component in the state space. Additionally, AoI-MDP introduces wait time in the action space, and integrate AoI with reward functions to achieve joint optimization of information freshness and decision making for AUVs leveraging RL for training. Simulation results highlight the feasibility, effectiveness and generalization of AoI-MDP over standard MDP, which effectively minimizes AoI while showcasing superior performance in the underwater task. The simulation code has been released as open-source to facilitate future research.

## REFERENCES

- [1] Z. Wang, Z. Zhang, J. Wang, C. Jiang, W. Wei, and Y. Ren, "Auv-assisted node repair for iout relying on multiagent reinforcement learning," *IEEE Internet of Things Journal*, vol. 11, no. 3, pp. 4139–4151, 2024.
- [2] Y. Li, L. Liu, W. Yu, Y. Wang, and X. Guan, "Noncooperative mobile target tracking using multiple auvs in anchor-free environments," *IEEE Internet of Things Journal*, vol. 7, no. 10, pp. 9819–9833, 2020.
- [3] R. H. Jhaveri, K. M. Rabie, Q. Xin, M. Chafii, T. A. Tran, and B. M. ElHalawany, "Guest editorial: Emerging trends and challenges in internet-of-underwater-things," *IEEE Internet of Things Magazine*, vol. 5, no. 4, pp. 8–9, 2022.
- [4] Z. Zhang, J. Xu, G. Xie, J. Wang, Z. Han, and Y. Ren, "Environment and energy-aware auv-assisted data collection for the internet of underwater things," *IEEE Internet of Things Journal*, vol. 11, no. 15, pp. 26406–26418, 2024.
- [5] W. Wei, J. Wang, J. Du, Z. Fang, C. Jiang, and Y. Ren, "Underwater differential game: Finite-time target hunting task with communication delay," in *ICC 2022 - IEEE International Conference on Communications*, 2022, pp. 3989–3994.
- [6] J. Wu, C. Song, J. Ma, J. Wu, and G. Han, "Reinforcement learning and particle swarm optimization supporting real-time rescue assignments for multiple autonomous underwater vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 7, pp. 6807–6820, 2022.
- [7] R. Talak, S. Karaman, and E. Modiano, "Optimizing information freshness in wireless networks under general interference constraints," *IEEE/ACM Transactions on Networking*, vol. 28, no. 1, pp. 15–28, 2020.
- [8] R. D. Yates, Y. Sun, D. R. Brown, S. K. Kaul, E. Modiano, and S. Ulukus, "Age of information: An introduction and survey," *IEEE Journal on Selected Areas in Communications*, vol. 39, no. 5, pp. 1183–1210, 2021.
- [9] K. Messaoudi, O. S. Oubbati, A. Rachedi, and T. Bendouma, "Uav-ugv-based system for aoi minimization in iot networks," in *ICC 2023 - IEEE International Conference on Communications*, 2023, pp. 4743–4748.
- [10] L. Lyu, Y. Dai, N. Cheng, S. Zhu, Z. Ding, and X. Guan, "Cooperative transmission for aoi-penalty aware state estimation in marine iot systems," in *2020 IEEE 18th International Conference on Industrial Informatics (INDIN)*, vol. 1, 2020, pp. 865–869.
- [11] Y. Sun, E. Uysal-Biyikoglu, R. D. Yates, C. E. Koksal, and N. B. Shroff, "Update or wait: How to keep your data fresh," *IEEE Transactions on Information Theory*, vol. 63, no. 11, pp. 7492–7508, 2017.
- [12] R. A. Howard, "Dynamic programming and markov processes," 1960. [Online]. Available: <https://api.semanticscholar.org/CorpusID:62124406>
- [13] E. Altman and P. Nain, "Closed-loop control with delayed information," *SIGMETRICS Perform. Eval. Rev.*, vol. 20, no. 1, p. 193–204, jun 1992.
- [14] B. Jiang, J. Du, C. Jiang, Z. Han, and M. Debbah, "Underwater searching and multiround data collection via auv swarms: An energy-efficient aoi-aware mappo approach," *IEEE Internet of Things Journal*, vol. 11, no. 7, pp. 12768–12782, 2024.
- [15] E. Altman and P. Nain, "Closed-loop control with delayed information," in *Proceedings of the 1992 ACM SIGMETRICS Joint International Conference on Measurement and Modeling of Computer Systems*, ser. SIGMETRICS '92/PERFORMANCE '92. New York, NY, USA: Association for Computing Machinery, 1992, p. 193–204.
- [16] K. Katsikopoulos and S. Engelbrecht, "Markov decision processes with delays and asynchronous cost collection," *IEEE Transactions on Automatic Control*, vol. 48, no. 4, pp. 568–574, 2003.