

# Adaptive AUV Hunting Policy with Covert Communication via Diffusion Model

Xu Guo\*, Xiangwang Hou<sup>†</sup>, Minrui Xu<sup>‡</sup>, Jianrui Chen<sup>§</sup>, Jingjing Wang<sup>§</sup>, Jun Du<sup>†</sup>, and Yong Ren<sup>†</sup>

\*Tsinghua Shenzhen International Graduate School, Tsinghua University, Shenzhen, 518055, China

<sup>†</sup>Department of Electronic Engineering, Tsinghua University, Beijing, 10084, China

<sup>‡</sup>School of Computer Science and Engineering, Nanyang Technological University, Singapore, 639798, Singapore

<sup>§</sup>School of Cyber Science and Technology, Beihang University, Beijing 100191, China

Email: guo-x24@mails.tsinghua.edu.cn, xiangwanghou@163.com, minrui001@e.ntu.edu.sg, chen\_jr@buaa.edu.cn, drwangjj@buaa.edu.cn, jundu@tsinghua.edu.cn, reny@tsinghua.edu.cn

**Abstract**—Collaborative underwater target hunting, facilitated by multiple autonomous underwater vehicles (AUVs), plays a significant role in various domains, especially military missions. Existing research predominantly focuses on designing efficient and high-success-rate hunting policy, particularly addressing the target’s evasion capabilities. However, in real-world scenarios, the target can not only adjust its evasion policy based on its observations and predictions but also possess eavesdropping capabilities. If communication among hunter AUVs, such as hunting policy exchanges, is intercepted by the target, it can adapt its escape policy accordingly, significantly reducing the success rate of the hunting mission. To address this challenge, we propose a covert communication-guaranteed collaborative target hunting framework, which ensures efficient hunting in complex underwater environments while defending against the target’s eavesdropping. To the best of our knowledge, this is the first study to incorporate the confidentiality of inter-agent communication into the design of target hunting policy. Furthermore, given the complexity of coordinating multiple AUVs in dynamic and unpredictable environments, we propose an adaptive multi-agent diffusion policy (AMADP), which incorporates the strong generative ability of diffusion models into the multi-agent reinforcement learning (MARL) algorithm. Experimental results demonstrate that AMADP achieves faster convergence and higher hunting success rates while maintaining covertness constraints.

**Index Terms**—Autonomous underwater vehicle (AUV), collaborative target hunting, covert communication, multi-agent reinforcement learning (MARL), diffusion model.

## I. INTRODUCTION

Multiple autonomous underwater vehicles (AUVs) enabled collaborative target hunting has been widely used in various fields, particularly military missions. However, collaborative underwater hunting tasks present numerous challenges including tracking, obstacle avoidance, and formation control [1]. This complexity has garnered significant academic attention, prompting extensive research in the field.

At the early stages, some works [2] [3] explore the problem of capturing stationary or slowly moving targets. These studies typically assume that the targets lack the ability to obtain information on the hunters’ positions or velocities. However, these assumptions are highly idealized. Real-world AUVs often possess advanced intelligence, enabling them to detect the hunters’ positions using sonar and other sensors, and to respond with evasion policy accordingly. As a further

development, the studies [4] and [5] consider a more practical scenario, where they assume that in the pursuit-evasion interaction between hunter AUVs and the target, the target could observe hunters’ locations. In fact, targets not only have observational abilities but also possess eavesdropping capabilities in actual scenarios. If sensitive information (e.g., collaborative policy, locations) exchanged within the formation of hunter AUVs is intercepted by the targets, they will adjust their evasion policy accordingly, which significantly impacts the success rate of the hunting process. However, existing research overlooks the impact of information leakage on the hunting process. To address this, we introduce covert communication techniques into the hunting framework. Covert communication [6] involves creating uncertainty in transmissions to prevent adversaries from intercepting sensitive information. Accordingly, we propose a covert communication-guaranteed collaborative target hunting framework that facilitates effective coordination among hunter AUVs while adhering to covert communication constraints, thereby preventing hunting policy from being compromised.

Moreover, coordinating multiple AUVs is an exceedingly complex task. Traditional rule-based methods for AUV target hunting [7] [8] require extensive parameter tuning to adapt to varying underwater conditions, yet they generally lack robustness and adaptability across different scenarios. In contrast, deep reinforcement learning (DRL) has proven to be an effective solution for target hunting [9] [10], boasting strong autonomous exploration capabilities. By continuously interacting with the environment, DRL enhances AUVs’ behavior through iterative learning and adaptation. However, DRL approaches often overlook modeling interactions and coordination among agents, which limits their effectiveness in collaborative tasks like multi-AUV hunting. To address these limitations, recent studies such as [4] and [11] introduce multi-agent reinforcement learning (MARL) for collaborative AUV hunting, leveraging MARL’s potential to optimize the joint trajectories of hunter AUV formations. Despite these advances, existing MARL frameworks for target hunting primarily rely on online RL, resulting in low data utilization. In contrast, offline RL offers improved data efficiency, which trains using pre-collected datasets, addresses this limitation. Therefore, we

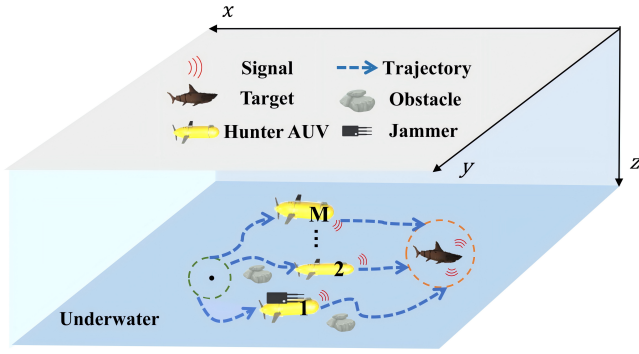


Fig. 1. Illustration of the AUVs hunting scenario with covert communication.

propose an adaptive multi-agent diffusion policy (AMADP) algorithm. AMADP utilizes the powerful policy generation capabilities of diffusion models to model the trajectories of hunter AUVs and integrates an adaptive attention mechanism to dynamically adjust formations, improving coordination among hunter AUVs while maintaining covert communication constraints. In summary, our contributions are as follows:

- As far as we know, this is the first attempt to consider the eavesdropping capabilities of the target and propose a covert communication-guaranteed collaborative target hunting framework, enabling efficient coordination under covert communication constraints.
- We design AMADP, a novel offline MARL algorithm, which leverages diffusion models to model hunter AUV trajectories and utilizes adaptive attention to adjust hunter formation under complex constraints.
- Extensive experimental results demonstrate that the proposed AMADP algorithm outperforms current state-of-the-art MARL algorithms in terms of hunting success rate and convergence speed under covert communication constraints.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

The system model is illustrated in Fig. 1. This scenario involves multiple hunter AUVs operating on a two-dimensional plane at a constant depth  $z$  underwater. The positions of hunter AUV  $i$  and the target at time  $t$  are represented by  $\mathbf{L}_i(t) = (x_i(t), y_i(t), z)$  and  $\mathbf{L}_T(t) = (x_T(t), y_T(t), z)$ , respectively. The hunter formation comprises  $M$  AUVs that collaboratively work to encircle the target while avoiding underwater obstacles. Hunter AUVs communicate with each other using a fixed transmission power  $P_S$ , exchanging critical information necessary for coordination. Furthermore, the first hunter AUV in the formation is equipped with a jammer that emits noise with power  $N_j$  to prevent the target from eavesdropping, ensuring covert communication throughout the hunting process.

### A. AUV Dynamics Model

In the hunting scenarios considered, the dynamic models of AUVs can be expressed by a simplified three-degree-of-

freedom model which describes the motion of the AUVs in the horizontal plane with a body-fixed coordinate frame  $\mathbf{v}_i = [w_i, v_i, r_i]^T$  and an earth-fixed reference frame  $\boldsymbol{\eta}_i = [u_{xi}, u_{yi}, \psi_i]^T$ , where  $w_i$ ,  $v_i$ ,  $r_i$ , and  $\psi_i$  represent the surge, sway, heave velocities, and yaw angle.  $\mathbf{v}_i$  is constrained by the maximum limit  $V_1$  satisfying  $\|\mathbf{v}_i\| \leq V_1$ . Similarly, we assume the velocity of the target  $\mathbf{v}_T$  is bounded by  $V_2$ , i.e.  $\|\mathbf{v}_T\| \leq V_2$ . The hunter AUVs and target use the same dynamics model [4], which be given by

$$\begin{cases} \dot{\boldsymbol{\eta}} = \mathbf{J}(\boldsymbol{\eta})\mathbf{v}, \\ \mathbf{M}_A\dot{\mathbf{v}} + \mathbf{C}_A(\mathbf{v})\mathbf{v} + \mathbf{B}_A(\mathbf{v})\mathbf{v} + \mathbf{G}_A(\boldsymbol{\eta}) = \mathbf{p} + \mathbf{e}, \end{cases} \quad (1)$$

where  $\mathbf{M}_A$ ,  $\mathbf{C}_A(\mathbf{v})$  and  $\mathbf{B}_A(\mathbf{v})$  represent the inertia matrix including added mass, the Coriolis-centripetal force matrix, and the damping matrix of the AUV, respectively. Additionally,  $\mathbf{G}_A(\boldsymbol{\eta})$  denotes the combined gravity and buoyancy matrix. The control input and environmental disturbance are represented by  $\mathbf{p}$  and  $\mathbf{e}$ . Besides,  $\mathbf{J}(\boldsymbol{\eta})$  is the transformation matrix, which can be given by

$$\mathbf{J}(\boldsymbol{\eta}) = \begin{bmatrix} \cos \psi & -\sin \psi & 0 \\ \sin \psi & \cos \psi & 0 \\ 0 & 0 & 1 \end{bmatrix}. \quad (2)$$

In the underwater target hunting task, hunter AUVs show better chasing ability by fulfilling cooperation, thus we assume the acceleration of hunter and the target satisfies  $\|\dot{\mathbf{v}}_i\| < \|\dot{\mathbf{v}}_T\|$ .

### B. Underwater Acoustic Channel Model

We formulate the path loss in shallow water environments as

$$A(d, f) = d^m a(f)^d, \quad (3)$$

where  $m$  represents the spreading factor characterizing the propagation geometry,  $f$  is the communication frequency, and  $d$  is the distance.  $a(f)$  denotes the attenuation coefficient per kilometer at a frequency  $f$  in kHz, is given by Thorp's empirical formula as

$$10 \lg a(f) = 3.3 \times 10^{-3} + \frac{0.11f^2}{1+f^2} + \frac{44f^2}{4100+f^2} + 3.0 \times 10^{-4}f^2. \quad (4)$$

Underwater environmental noise includes turbulence noise  $N_t$ , shipping noise  $N_s$ , wind-driven wave noise  $N_w$ , and thermal noise  $N_{th}$ . Based on the underwater noise model, the total noise level can be expressed as

$$\begin{cases} 10 \lg N_t(f) = 17 - 30 \lg f, \\ 10 \lg N_s(f) = 40 + 20(s - 0.5) + 26 \lg f - 60 \lg(f + 0.03), \\ 10 \lg N_w(f) = 50 + 7.5\sqrt{w} + 20 \lg f - 40 \lg(f + 0.4), \\ 10 \lg N_{th}(f) = -15 + 20 \lg f, \end{cases} \quad (5)$$

where  $s$  and  $w$  denote shipping activity factor and wind speed. These noise sources are modeled as Gaussian processes, and the total underwater environment noise power spectral density is given by  $N_u(f) = N_t(f) + N_s(f) + N_w(f) + N_{th}(f)$  in dB re  $\mu\text{Pa}$  per Hz. The acoustic path loss and underwater noise power at  $t$ -th time step are denoted as  $A[t]$  and  $N_u[t]$ .

### C. Covert Communication Model

We assume that the channel state is constant during each time step but changes independently from one time step to the next. To describe the underwater time-varying channel state, we use a block fading channel model, considering each channel to be mutually independent. The hunter AUVs transmit signals through multiple channels during communication, represented as  $\mathbf{y}^R[t] = [y_1^R[t], \dots, y_L^R[t]]$ , where  $L$  denotes the number of channels uses. Meanwhile, the target passively monitors the signal strength to detect communication among hunter AUVs. The signal received by the target on the  $l$ -th channel in the  $t$ -th time step is expressed as [6]

$$y_l^T[t] = \begin{cases} \sqrt{\frac{P_S[t]}{A_{ac}[t]}} s_l + n_l^T[t], & H_1, \\ n_l^T[t], & H_0, \end{cases} \quad (6)$$

where hypothesis  $H_1$  represents that the hunter AUVs are communicating with each other, and  $H_0$  indicates they are not.  $P_S[t]$  and  $s_l$  denote the transmit power and the transmitted signal on the  $l$ -th channel in the  $t$ -th time step, respectively. Besides,  $n_l^T[t] \sim \mathcal{N}(0, N_T[t])$ , where  $N_T[t] = N_u[t] + N_j[t]$ , represents the total Gaussian noise observed by the target in the  $t$ -th time step and  $A_{ac}[t]$  refers to the acoustic path loss experienced by the signal as it propagates from the hunter AUVs to the target.

In this scenario, the target infers the communication state between hunter AUVs by measuring the received signal power and applying a hypothesis testing strategy to aid its decision-making during the hunting task. Specifically, to detect the presence of the covert communication, target needs to clarify whether a hunter AUV sends information to others. During the hunting process, the target employs the likelihood ratio test (LRT) [12] as the optimal detection approach to determine whether the hunter AUVs are communicating, while minimizing its detection error. This process can be simplified as

$$Y[t] \triangleq \frac{1}{L} \sum_{l=1}^L |\tilde{y}_l^T[t]|^2 \underset{H_0}{\overset{H_1}{\geq}} \alpha[t], \quad (7)$$

where  $Y[t]$  represents the average received signal power of the target,  $\sum_{l=1}^L |\tilde{y}_l^T[t]|^2$  is the total received power at the target over the block duration, and  $\alpha[t]$  is the threshold for the  $t$ -th time slot. If the received power is below the threshold, the target favors hypothesis  $H_0$ ; otherwise,  $H_1$ .

Under the assumption that target knows the transmission power  $P_S[t]$  of the hunter AUVs and noise power  $N_T[t]$ , and the optimal threshold  $\alpha^*[t]$  at target is given by [13]

$$\alpha^*[t] = N_T[t] \left( 1 + \frac{1}{\beta_T[t]} \right) \ln(1 + \beta_T[t]), \quad (8)$$

where  $\beta_T[t]$  equals  $\frac{P_S[t]}{A_{ac}[t]N_T[t]}$ .

Let  $P_{FA}$  represent the probability of false alarm, indicating that the target favors hypothesis  $H_1$  when the actual state is

$H_0$ . Similarly,  $P_{MD}$  denote the missed detection probability indicating a preference for  $H_0$  when  $H_1$  is true. To ensure covert communication, the following constraint must be satisfied

$$P_{FA} + P_{MD} \geq 1 - \epsilon, \quad (9)$$

where  $\epsilon$  signifies the acceptable level of covertness. However, calculating  $P_{FA}[t]$  and  $P_{MD}[t]$  directly is challenging. Therefore, according to [13], the constraint equation can be transformed into

$$\frac{L}{2} \left[ \ln(1 + \beta_T[t]) - \frac{\beta_T[t]}{1 + \beta_T[t]} \right] < 2\epsilon^2. \quad (10)$$

### D. Problem Formulation

The proposed covert communication-guaranteed collaborative target hunting framework requires completing the hunting successfully while ensuring covert communication constraints throughout the entire process.

**Successful hunting criteria:** We assume that the detection range and the attacking range of the hunter AUVs are  $R_1$  and  $R_2$ , respectively. Specifically, when the distance between the target and the  $i$ -th hunter AUV  $e_i$  is less than  $R_1$  ( $\|e_i\| < R_1$ ), the hunter AUVs obtain the target's location information and share it within the formation. hunting is successful if all  $M$  hunter AUVs are within the distance  $R_2$  from the target ( $\|e_i\| < R_2$ ) and form an encirclement. Conversely, if at the end of the given time steps  $h$  all  $\|e_i\| > R_1$ , or if the conditions for a successful hunting are not met within this time, the hunting operation is considered a failure.

**Covert communication constraint:** The total divergence between the probability distributions under hypotheses  $H_0$  and  $H_1$  is constrained as

$$D_{KL}(Q_0 \| Q_1)[t] = \frac{L}{2} \left[ \ln(1 + \beta_T[t]) - \frac{\beta_T[t]}{1 + \beta_T[t]} \right] \leq 2\epsilon^2, \quad (11)$$

where  $\epsilon$  represents the error tolerance, and  $\mathcal{D}_{KL}(Q_0 \| Q_1)[t]$  is the Kullback-Leibler (KL) divergence between the distributions under different hypotheses at time  $t$ .  $Q_1$  and  $Q_0$  represents the probability density under hypothesis  $H_1$  and  $H_0$ , respectively.

**Remark 1:** For simplicity, we assume that the transmission power  $P_S[t]$  and communication frequency  $f$  between hunter AUVs and the jammer power  $N_j[t]$  directed at the target all remain constant. Therefore the covert communication constraints are solely dependent on the distances between the hunters and the target. Our goal is to optimize the trajectories of the hunter AUVs in the formation to complete the hunting task while satisfying the covert communication constraints.

In summary, we define the optimization problem as

$$\begin{aligned} \max P_{\text{success}} &= \Pr(\|e_i\| < R_2 \quad \forall i = 1, \dots, M) \\ \text{s.t.} & \begin{cases} D_{KL}(Q_0 \| Q_1)[t] \leq 2\epsilon^2, & \forall t, \\ \|\mathbf{v}_i(t)\| \leq V_1, \quad \|\mathbf{v}_T(t)\| \leq V_2, & \forall t, \\ \|\mathbf{L}_i(t) - \mathbf{L}_j(t)\| \geq r_{\min}, & \forall i \neq j. \end{cases} \end{aligned} \quad (12)$$

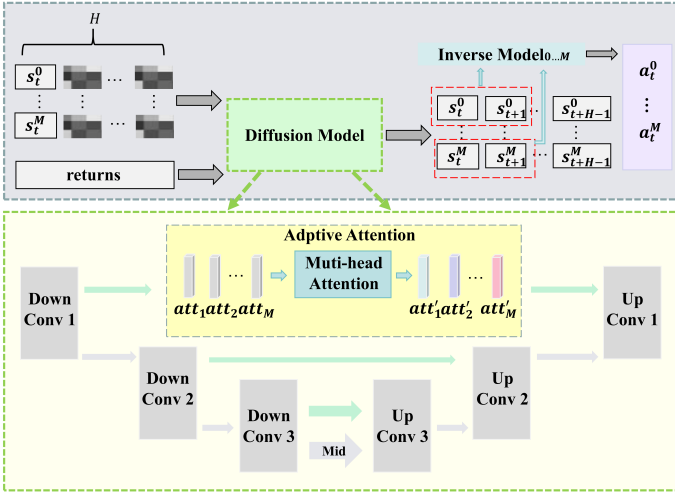


Fig. 2. Illustration of the AUVs hunting scenario with covert communication.

where  $P_{\text{success}}$  represents the probability of successful encirclement when all hunter AUVs are within the target's capture range  $R_2$ .  $D_{\text{KL}}(Q_0||Q_1)[t]$  enforces covert communication constraints,  $V_1$  and  $V_2$  are the maximum speeds for the hunters and the target, respectively, and  $r_{\text{min}}$  is the minimum safe distance to avoid collisions.

### III. METHODOLOGY

#### A. Markov Game Process Modeling

We model the hunting process as a partially observable Markov decision process (POMDP), where each AUV and the target act as agents

$$M = (S_i, A_i, P_i, R_i), \quad (13)$$

where  $S_i$  represents the observation information of each agent,  $A_i$  is the action taken by each agent,  $P_i(s'|s, a)$  determines the transition probability to the next state, and  $R_i$  is the reward information. In this context, hunter AUVs receive the same rewards.

1) *State Space*: The observation information for each agent can be defined as

$$S_i(t) = \left\{ v_i(t), \mathbf{L}_i(t), \{\mathbf{L}_o(t)\}, \{\mathbf{L}_j(t)\}_{j \neq i} \right\}, \quad (14)$$

where  $v_i(t)$  and  $\mathbf{L}_i(t)$  represent the velocity and position information of agent  $i$ , respectively.  $\{\mathbf{L}_o(t)\}$  denotes the positions of the obstacles, and  $\{\mathbf{L}_j(t)\}_{j \neq i}$  represents the positions of other agents.

2) *Action Space*: Each agent's action  $a_i(t)$  includes the movement direction  $\theta_i$  and velocity  $v_i$ , which is represented as

$$a_i(t) = \{\theta_i(t), v_i(t)\}. \quad (15)$$

3) *Reward Function*: To achieve effective encirclement and covert communication constraints under dynamic and uncertain underwater conditions, the reward function is designed to encourage the hunter AUVs to form an optimal encirclement around the target while maintaining specific communication constraints.

- **Encirclement formation reward**  $R_i^E(t)$ : This reward encourages hunters to form a well-distributed encirclement around the target,

$$R_i^E(t) = \begin{cases} -\lambda\sigma_s(t), & d_g(t) > d_g^*, \\ \zeta(d_g^* - d_g(t)) - \lambda\sigma_s(t), & d_g(t) \leq d_g^*, \end{cases} \quad (16)$$

where  $\sigma_s(t)$  represents the variance in the distances between the hunters to encourage a balanced formation, and  $d_g(t)$  is the distance between the target and the centroid of the hunter AUVs formation at time  $t$ , and  $d_g^*$  is the desired distance for encirclement completion. When  $d_g(t) \leq d_g^*$ , an additional reward is provided for a successful hunting task.

- **Collision avoidance reward**  $R_i^C(t)$ : This reward penalizes collisions between hunter AUVs or with landmarks, helping maintain a safe distance during the hunting task.

$$R_i^C(t) = \begin{cases} -\nu, & \text{if collision,} \\ 0, & \text{otherwise.} \end{cases} \quad (17)$$

- **Covert performance reward**  $R_i^P(t)$ : This reward is designed to ensure that the hunting process adhere to the covert communication constraint  $\epsilon$ .

$$R_i^P(t) = \begin{cases} \nu, & D(Q_0||Q_1)[t] \leq 2\epsilon^2, \\ -\nu, & D(Q_0||Q_1)[t] > 2\epsilon^2. \end{cases} \quad (18)$$

The overall reward for an hunter AUV  $R_i(t)$  is a sum of the above rewards

$$R_i(t) = R_i^E(t) + R_i^C(t) + R_i^P(t), \quad (19)$$

This reward design ensures that hunters effectively perform hunting tasks while avoiding collisions and meeting covert communication requirements.

#### B. Design of AMADP Algorithm

1) *Diffusion Model*: The proposed algorithm is based on the denoising diffusion probabilistic model (DDPM) [14], a powerful deep generative model designed to learn the underlying data distribution  $q(\mathbf{x}_0)$  from a dataset  $\mathcal{D} = \{\mathbf{x}_k\}$ . In DDPM, data reconstruction is achieved by denoising real data  $\mathbf{x}_0$  from noises  $\mathcal{N}(0, \mathbf{I})$  over  $K$  diffusion steps.

The predefined forward process is defined as  $q(\mathbf{x}_k|\mathbf{x}_{k-1}) = \mathcal{N}(\sqrt{\alpha_k}\mathbf{x}_{k-1}, \sqrt{1-\alpha_k}\mathbf{I})$ , where  $\alpha_k = 1 - \beta_k$ , and  $\beta_{1:K}$  is a variance schedule. The reverse process is parameterized as  $p_\theta(\mathbf{x}_{k-1}|\mathbf{x}_k) = \mathcal{N}(\mu_\theta(\mathbf{x}_k), \Sigma_k)$ . The mean  $\mu_\theta(\mathbf{x}_k)$  and variance  $\Sigma_k$  can be expressed as  $\mu_\theta(\mathbf{x}_k) = \frac{1}{\sqrt{\alpha_k}} \left( \mathbf{x}_k - \frac{\beta_k}{\sqrt{1-\alpha_k}} \epsilon_\theta(\mathbf{x}_k, k) \right)$  and  $\Sigma_k = \beta_k \frac{1-\bar{\alpha}_{k-1}}{1-\alpha_k} \mathbf{I}$ , where  $\bar{\alpha}_k = \prod_{i=1}^k \alpha_i$ .

The loss function of DDPM is defined as

$$L(\theta) = \mathbb{E}_{k \sim \{1, \dots, K\}, \mathbf{x}_0, \epsilon} [\|\epsilon - \epsilon_\theta(\mathbf{x}_k, k)\|^2], \quad (20)$$

where  $\epsilon$  is the real noise added in each diffusion step, and  $\epsilon_\theta(\mathbf{x}_k, k)$  is the noise predicted by the noise prediction network at diffusion step  $k$ .

2) *Architecture of AMADP*: The overall architecture of the network is illustrated in Fig. 2. The entire AUV formation shares a noise prediction network, which is based on a modified U-Net architecture consisting of three down-sampling convolutional layers and three up-sampling convolutional layers, connected by a middle bottleneck. The down-sampling layers progressively compress the state trajectories features while integrating the conditioning information. The up-sampling layers reconstruct the state trajectories by concatenating the intermediate features from the corresponding down-sampling layers via skip connections. After the modified U-Net predicts state trajectories for each hunter AUV. These trajectories are then fed into the respective inverse dynamics models, allowing each hunter to determine the necessary actions at each time step for optimal target hunting.

To coordinate the formation between multiple hunter AUVs, we introduce an adaptive attention mechanism between the initial down-sampling and the up-sampling layer of the U-Net. Unlike conventional self-attention mechanisms that focus solely on global information, the adaptive attention is designed to dynamically adjust attention weights based on both global and agent-specific inputs. Formally, adaptive attention can be expressed as

$$\begin{cases} att_i = \text{softmax}\left(\frac{Q_{i,t} \cdot K_t^T}{\sqrt{d_k}}\right) \cdot V_t, \\ att = \text{concat}(att_1, att_2, \dots, att_M), \end{cases} \quad (21)$$

where  $Q_{i,t}$  is the query of the  $i$ -th AUV,  $K_t$  and  $V_t$  represent the global key and value information.  $att_i$  represents the attention weight for each AUV. The attention weight  $att$  reflects the global attention information, incorporating the contributions of each hunter during the hunting process. By dynamically adjusting these attention weights, the model can better capture the interactions between hunter AUVs and optimize their coordination for effective target hunting.

3) *Training Framework of AMADP*: Unlike online RL, which requires real-time interaction with the environment for continuous policy updates, offline RL relies on a pre-collected static dataset  $\mathcal{D}$  to learn the policy, thereby improving data utilization. Considering the challenges posed by obstacles and unstable communication in underwater environments make real-time interaction extremely challenging, AMADP adopts an offline RL training approach. Moreover, diffusion-based offline RL is suitable for solving cooperative games. Thus the hunting policy of the entire hunter AUV formation is generated by the AMADP algorithm proposed in this paper, and the escape policy of the target is the pre-trained traditional RL method deep deterministic policy gradient (DDPG) [15].

In our approach, the diffusion model is conditioned on  $y_\tau$ , which contains the current state, achieved return, and the current time step to generate future trajectories. Considering the actual hunting process, AMADP uses a centralized training with decentralized execution (CTDE) framework, where global

TABLE I. Parameters of System and Algorithm

Parameters	Values
Diffusion step ( $K$ )	200
Learning rate	0.0001
Training step	20000
Batch size	32
Discounting factor	0.9
Return scale	3000
Planning ahead time steps ( $H$ )	40
Reward design ( $\lambda, \zeta, \nu$ )	20, 6000, 10
Start point of hunter AUVs ( $O$ )	(500, 500, -200) m
Number of hunter AUVs ( $M$ )	3
Maximum speed of hunter AUV ( $V_1$ )	0.3 m/s
Maximum speed of target ( $V_2$ )	0.2 m/s
Acceleration of hunter AUV ( $\ \dot{v}_i\ $ )	0.01 m/s <sup>2</sup>
Acceleration of target ( $\ \dot{v}_T\ $ )	0.02 m/s <sup>2</sup>
Movement range of AUV ( $\psi$ )	$[-\pi, \pi]$
Sensing radius of AUV ( $R_1$ )	800 m
Attacking radius of AUV ( $R_2$ )	150 m
Desired distance ( $d_g^*$ )	120 m
Communication Constraint ( $\epsilon$ )	0.04
Communication Frequency ( $f$ )	25 kHz
Transmission Power ( $P_S$ )	0.1 W
Jammer Power ( $N_T$ )	0.2 W

information is accessed during training, but each hunter AUV makes decisions based on local observations during execution. To simplify the representation of the decision process in the diffusion model, we define the learning state sequence as

$$\tau = [s_0^1, \dots, s_0^M, s_1^1, \dots, s_1^M, \dots, s_H^1, \dots, s_H^M], \quad (22)$$

where  $s_t^i$  indicates the state of the  $i$ -th hunter AUV at time step  $t$  and  $H$  represents the planning ahead time steps.

For each hunter AUV, we define an inverse dynamics model  $f_\phi$ , which predicts the action  $a_t^i$  of the  $i$ -th AUV at time step  $t$  based on its current state  $s_t^i$  and next state  $s_{t+1}^i$

$$a_t^i = f_\phi(s_t^i, s_{t+1}^i). \quad (23)$$

Incorporating the DDPM loss, the inverse dynamics model loss, and the classifier-free guidance mechanism [16], the overall training loss function can be formulated as

$$\begin{aligned} \mathcal{L}(\theta, \phi) := & \sum_i \mathbb{E}_{(s_t^i, a_t^i, s_{t+1}^i) \in \mathcal{D}} \left[ \|a_t^i - f_\phi^i(s_t^i, s_{t+1}^i)\|^2 \right] \\ & + \mathbb{E}_{k, \tau \in \mathcal{D}, \beta} \left[ \|\epsilon - \epsilon_\theta(\hat{\tau}_k, (1 - \beta)y(\hat{\tau}) + \beta\theta, k)\|^2 \right], \end{aligned} \quad (24)$$

where  $\beta$  is sampled from a Bernoulli distribution to balance the conditioned and unconditioned diffusion process.

#### IV. EXPERIMENTAL RESULTS

The experiments are conducted in a 1200 m  $\times$  1200 m area with a water depth of -200 m, where obstacles are randomly distributed. The AUV formation starts from the central point (500, 500) and the target's position is randomly initialized. The parameters of the system model and the AMADP algorithm are presented in Table I.

Fig. 3(a) illustrates a successful encirclement using the AMADP algorithm, where the AUV formation successfully surrounds the target while avoiding obstacles without any collisions. We measured the  $D_{\text{KL}}(Q_0 \| Q_1)$  values at each time step over 50 episodes, taking the average value of

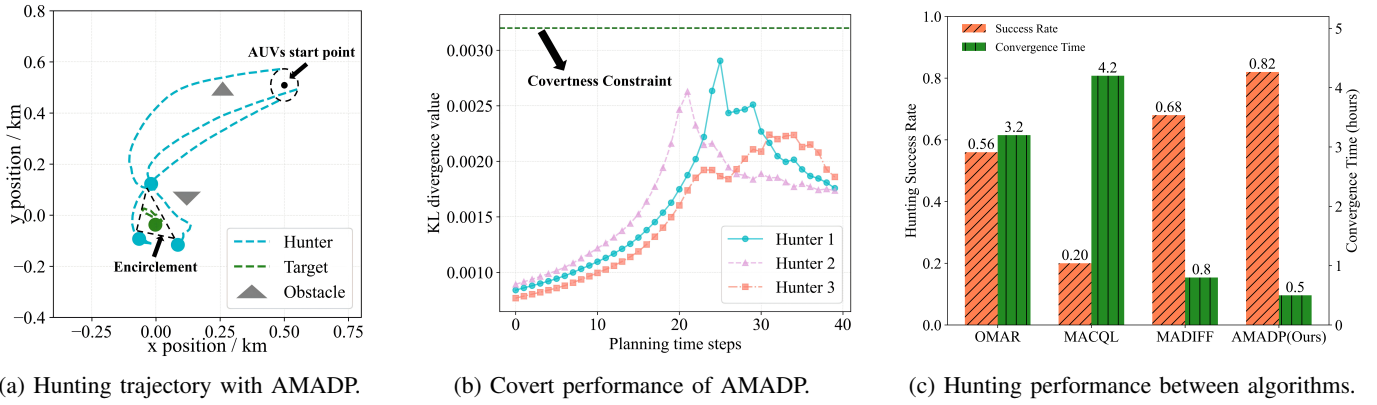


Fig. 3. Hunting with covert communication between hunter AUVs and the target.

each time step, as shown in Fig. 3(b). The results indicate that the AMADP algorithm consistently satisfies the covert communication constraints throughout the hunting process.

To demonstrate the superior performance of the proposed algorithm, we compare AMADP with state-of-the-art offline MARL algorithms, multi-agent conservative Q-learning (MACQL) [17], OMAR [18] and MADIFF [19]. The experimental results in Fig. 3(c) indicate that AMADP significantly outperforms other algorithms in the success rate of hunting task. And also due to the simplified structure of our proposed algorithm, it achieves faster convergence.

## CONCLUSION

In this paper, we consider the eavesdropping capabilities of the target during collaborative hunting and propose a covert communication-guaranteed hunting framework for the first time. To improve the generalization capability, data efficiency, and trajectory diversity of traditional collaborative hunting algorithms, we introduce AMADP, an offline MARL algorithm, which incorporates diffusion models to generate diverse hunting trajectories and utilizes adaptive attention to dynamically adjust the formation. Experimental results demonstrate that AMADP satisfies covert communication constraints and achieves higher success rates and faster convergence compared to existing state-of-the-art algorithms.

## REFERENCES

- [1] M. Zhang, H. Chen, and W. Cai, "Hunting task allocation for heterogeneous multi-auv formation target hunting in iout: A game theoretic approach," *IEEE Internet of Things Journal*, vol. 11, no. 5, pp. 9142–9152, 2024.
- [2] C. Wang and G. Xie, "Limit-cycle-based decoupled design of circle formation control with collision avoidance for anonymous agents in a plane," *IEEE Transactions on Automatic Control*, vol. 62, no. 12, pp. 6560–6567, 2017.
- [3] M. Deghat, I. Shames, B. D. Anderson, and C. Yu, "Localization and circumnavigation of a slowly moving target using bearing measurements," *IEEE Transactions on Automatic Control*, vol. 59, no. 8, pp. 2182–2188, 2014.
- [4] W. Wei, J. Wang, J. Du, Z. Fang, Y. Ren, and C. L. P. Chen, "Differential game-based deep reinforcement learning in underwater target hunting task," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–13, 2023.
- [5] C. Wang, Y. Wang, P. Shi, and F. Wang, "Scalable-maddpg-based cooperative target invasion for a multi-usv system," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–11, 2023.
- [6] X. Chen, J. An, Z. Xiong, C. Xing, N. Zhao, F. R. Yu, and A. Nalnanathan, "Covert communications: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 25, no. 2, pp. 1173–1198, 2023.
- [7] X. Meng, B. Sun, and D. Zhu, "Harbour protection: moving invasion target interception for multi-auv based on prediction planning interception method," *Ocean Engineering*, vol. 219, p. 108268, 2021.
- [8] C. Lin, G. Han, J. Du, Y. Bi, L. Shu, and K. Fan, "A path planning scheme for auv flock-based internet-of-underwater-things systems to enable transparent and smart ocean," *IEEE Internet of Things Journal*, vol. 7, no. 10, pp. 9760–9772, 2020.
- [9] T. Ma, J. Lyu, J. Yang, R. Xi, Y. Li, J. An, and C. Li, "Clsql: Improved q-learning algorithm based on continuous local search policy for mobile robot path planning," *Sensors*, vol. 22, no. 15, 2022.
- [10] M. Xi, J. Yang, J. Wen, H. Liu, Y. Li, and H. H. Song, "Comprehensive ocean information-enabled auv path planning via reinforcement learning," *IEEE Internet of Things Journal*, vol. 9, no. 18, pp. 17440–17451, 2022.
- [11] Z. Wang, Y. Sui, H. Qin, and H. Lu, "State super sampling soft actor-critic algorithm for multi-auv hunting in 3d underwater environment," *Journal of Marine Science and Engineering*, vol. 11, no. 7, 2023.
- [12] M. K. Steven, "Fundamentals of statistical signal processing," *PTR Prentice-Hall, Englewood Cliffs, NJ*, vol. 10, no. 151045, p. 148, 1993.
- [13] J. Chen, J. Wang, Z. Wei, Y. Ren, C. Masouros, and Z. Han, "Joint autonomous underwater vehicle trajectory and energy optimization for underwater covert communications," *IEEE Transactions on Communications*, pp. 1–1, 2024.
- [14] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [15] T. Lillicrap, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [16] J. Ho and T. Salimans, "Classifier-free diffusion guidance," *arXiv preprint arXiv:2207.12598*, 2022.
- [17] A. Kumar, A. Zhou, G. Tucker, and S. Levine, "Conservative q-learning for offline reinforcement learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 1179–1191, 2020.
- [18] L. Pan, L. Huang, T. Ma, and H. Xu, "Plan better amid conservatism: Offline multi-agent reinforcement learning with actor rectification," in *International conference on machine learning*. PMLR, 2022, pp. 17221–17237.
- [19] Z. Zhu, M. Liu, L. Mao, B. Kang, M. Xu, Y. Yu, S. Ermon, and W. Zhang, "Madiff: Offline multi-agent learning with diffusion models," *arXiv preprint arXiv:2305.17330*, 2023.