

Cooperative guidance of multiple missiles: a hybrid co-evolutionary approach

Xuejing Lan, Junda Chen, Zhijia Zhao, *Member, IEEE*, and Tao Zou

Abstract—Cooperative guidance of multiple missiles is a challenging task with rigorous constraints of time and space consensus, especially when attacking dynamic targets. In this paper, the cooperative guidance task is described as a distributed multi-objective cooperative optimization problem. To address the issues of non-stationarity and continuous control faced by cooperative guidance, the natural evolutionary strategy (NES) is improved along with an elitist adaptive learning technique to develop a novel natural co-evolutionary strategy (NCES). The gradients of the original evolutionary strategy are rescaled to reduce the estimation bias caused by the interaction between the multiple missiles. A hybrid co-evolutionary cooperative guidance law (HCCGL) is then developed by integrating the highly scalable co-evolutionary strategy and the proportional guidance law, with detailed convergence proof provided. Finally, simulations demonstrated the effectiveness and superiority of this guidance law in solving cooperative guidance tasks with high accuracy, with potential applications in other multi-objective optimization, dynamic optimization, and distributed control scenarios.

Index Terms—Optimal control; cooperative guidance; evolutionary strategy; multi-objective optimization

I. INTRODUCTION

MODERN penetration of air defense systems of the target requires coordinated attacks with multiple missiles. However, the rapid development of detection technologies and close-in weapon systems (CIWS) has decreased the chances of successful impact with a single conventional missile. [1]. In addition to increasing the difficulty of interception, the cooperative guidance strategy of multiple missiles is also crucial to the lethal effect of the final impact. Usually, the cooperative guidance of multiple missiles belongs to the phase of terminal guidance, where accurate target information can be obtained with active radar systems or other detection devices. The existing cooperative guidance laws can be roughly divided into two categories. One is the analytical method to find closed-form solutions, which is mainly based on sliding mode control, optimal control, and multi-agent consensus theory. The other is the intelligent method which generally adopts heuristic intelligent optimization algorithm and reinforcement learning (RL) theory.

The analytical cooperative guidance method has been proven to be robust and efficient for practical application [2]–[6]. Based on fundamental proportional navigation (PN), Jeon et.

al developed cooperative proportional navigation (CPN) where the on-board time-to-go of the missile is used as the navigation gain [1]. It is a simple but effective approach for achieving time consensus. Ma developed a composite guidance law, which can be decomposed into the direction along the line of sight (LOS) and the direction perpendicular to LOS [2], corresponding to time and space cooperative respectively. Furthermore, time cooperative control is achieved with the combination of PNG and impact time error feedback [7], where the undirected topology is adopted to establish communication relationships. Based on the optimal control approach, a variant of the hyperbolic tangent function is proposed in [3] to force early control of velocity and impact angle.

However, with the increasing demand for developing high-precision weapon systems, intelligent cooperative guidance method is increasingly regarded as a necessary auxiliary option. In recent years, the reinforcement learning theory has attracted much attention because of its ability to learn online based on environmental feedback [8]–[13]. According to the training structures, existing reinforcement learning algorithms for multi-agent systems can be roughly divided into four types, which are *Fully decentralized training, decentralized execution; Fully centralized training, decentralized execution; Centralized training, centralized execution, and value decomposition methods*. Some of these algorithms have achieved satisfactory results in coping with problems with low complexity and accuracy requirements. In [14], [5], and [15], the state-of-the-art reinforcement learning frameworks have demonstrated their effectiveness in the guidance task. Zhang et.al proposed a gradient-descent-based reinforcement learning method in the actor-critic framework and achieved consensus control for multi-agent systems by following a tracking leader [16]. But the two challenges of *Nonstationarity* and *Partial Observability* [17] will lead to saturated output or coordination loss of multi-agent systems, which greatly reduces the accuracy of the value function. In addition, the use of value function in reinforcement learning is not suitable for continuous control tasks with large search spaces. Thus, these limitations of RL impede the development of reinforcement learning in cooperative guidance.

It is an excellent way to solve the above problems by removing the value function of reinforcement learning and optimizing in solution space with evolutionary strategy (ES), which is more robust and invariant to real-time rewards because it optimizes towards the objective function directly [18]. Moreover, as described in [19], ES is tolerant of long horizontal and implicit solutions, which is exactly consistent with the need for cooperative guidance. The natural evolutionary strategy (NES) is the latest branch of ES, and shows good performance in

The author(s) received no financial support for the research, authorship, and/or publication of this article. (Xuejing Lan and Junda Chen contributed equally to this work)(Corresponding author: Zhijia Zhao)

X. J. Lan, J. D. Chen, Z. J. Zhao, and T. Zou are with the School of Mechanical and Electrical Engineering, Guangzhou University, Guangzhou 510006, China (e-mail: lanxj@gzhu.edu.cn; CJD@e.gzhu.edu.cn; zhjzhaos-cut@163.com; tzou@gzhu.edu.cn).

solving high-dimensional continuous multimodal optimization problems, by using the natural gradient information estimated according to the fitness expectation of the population [18]–[20]. Similar algorithms named co-evolutionary algorithm have been discussed in [21] and [22], which focus on solving multi-objective optimization problems by dividing the overall objective into sub-objectives, such to optimize and evaluate together. Another idea is to evolve multiple populations for the same goal, and manually regulate the constraints of each population for faster convergence or fuller exploration [22]. As represented in [22], the concept of co-evolution refers to multi-threads of training processes. Note that these methods do not use the natural gradient information as in NES, and the non-stationary issue discussed above is not considered.

When optimizing in continuous parameter(solution) space, it is very important to apply adaptive technology. While a learning rate adaption method based on the quality of gradients is often not easy to estimate, a simple workaround would be leveraging the shifting distance of parameters to adapt the learning rate. As shown in [23], the size of population was adjusted depending on the novelty metric and quantity metric, which reflected the complexity of the dynamic environment. The estimation of distribution algorithm (EDA) was applied to continuous control by searching the optimal parameter distribution [24], [25]. A variety of evolutionary methods were investigated to design the multi-objective missile guidance law [26]. Maheswaranathan proposed a surrogate gradient to reduce the evaluation costs [27]. These works reveal the enormous potential of searching in parameter space, rather than directly searching in parameter space.

Therefore, an NES-based co-evolutionary algorithm naming as the natural co-evolutionary strategy (NCES) is developed in this paper to distress the dilemma faced by RL in the cooperative guidance task. Considering the advantages of searching in parameter space, the co-evolutionary algorithm is improved in this work by rescaling the gradient information to reduce the estimation bias introduced by neighboring populations. As discussed in [28], [29], most of today's bio-inspired algorithm innovations are based on experimental observation rather than meticulous theoretical support. Whereas in this work, we try to dig into the depths of complex optimization and provide proof as sensible as possible through the presentation of graphs and deduction. Via integrating the NCES algorithm, a hybrid co-evolutionary cooperative guidance law (HCCGL) is further developed to solve the challenging missile guidance problem. Extensive empirical results on various engagement scenarios verified the effectiveness of the proposed guidance law. The main contributions of this work are summarized as follows:

- 1) To address the issues of non-stationarity and continuous control faced by cooperative guidance, an NCES algorithm is formulated and incorporated into a novel guidance law as an alternative to RL in the cooperative guidance task.
- 2) The rigorous constraints of time and space consensus in cooperative guidance are integrated and designed as the fitness function for each missile. An MLP-based policy network is constructed and learned to optimize the fitness function.

- 3) The proposed HCCGL has advantages in achieving high precision for cooperative guidance tasks, even with dynamic targets and random initial conditions.

The rest of the paper is organized as follows. The problem formulation is elaborated in Section II, and the proposed cooperative guidance law is discussed in Section III. In Section V, experiments under various configurations are implemented. Finally, conclusions are made in Section VI.

II. PROBLEM FORMULATION

A. Engagement geometry

The two-dimensional engagement geometry between multiple missiles and one target is shown in Fig. 1, where the inertial coordinate frame OXY represents the horizontal plane. There are n missiles in total. The index M_i denotes the i^{th} missile, and T represents the target. V_{mi} , Ξ_{mi} , α_{mi} , and δ_{mi} represent the velocity, line of sight (LOS) angle, flight-path angle, and heading angle of the i^{th} missile, respectively. a_{li} and a_{vi} represent the lateral acceleration and the thrust acceleration to be designed for the i^{th} missile, which are perpendicular to and align with the direction of V_{mi} , respectively. V_T , Ξ_T , α_T , and δ_T are the velocity, LOS angle, flight-path angle, and heading angle of the target, respectively. The lateral acceleration of the target is denoted by a_T .

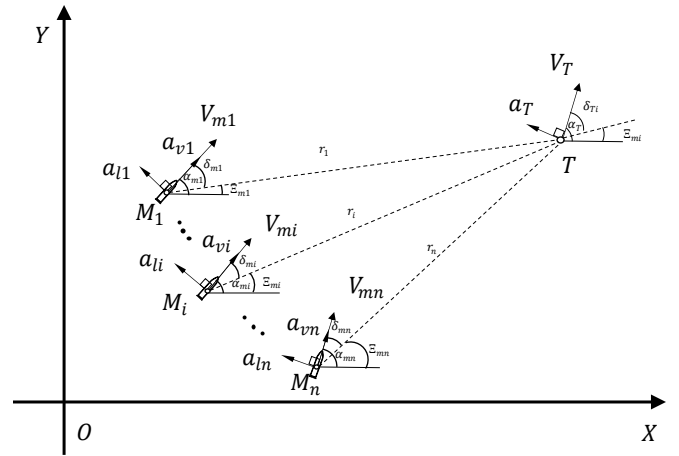


Fig. 1: Two-dimensional engagement geometry

The dynamic equations of the i^{th} missile and the target are as follows:

$$\begin{cases} \dot{r}_i = -V_{mi}\cos\delta_{mi} + V_T\cos\delta_{Ti} \\ r_i\dot{\Xi}_{mi} = -V_{mi}\sin\delta_{mi} + V_T\sin\delta_{Ti} \\ \dot{\alpha}_{mi} = a_{li}/V_{mi} \\ \dot{\alpha}_T = a_T/V_T \\ \dot{V}_{mi} = a_{vi} \\ \delta_{mi} = \alpha_{mi} - \Xi_{mi} \\ \delta_{Ti} = \alpha_T - \Xi_{mi} \end{cases}, \quad (1)$$

where, r_i represents the relative range between the i^{th} missile and the target. The time-to-go of the i^{th} missile t_{go}^i refers to the time left from the current time until the interception:

$$t_{go}^i = -\frac{r_i}{\dot{r}_i}, \quad (2)$$

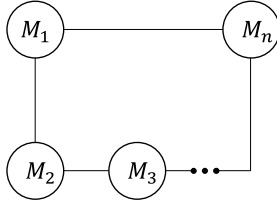


Fig. 2: Communication topology

B. Communication Topology

The communication relationship of the multiple missiles is depicted by a topology, where a set of nodes $\mathcal{V} = \{v_1, v_2, \dots, v_n\}$ represents the n missiles. The communications are represented by a set of edges $\xi \subseteq \mathcal{V} \times \mathcal{V}$ with an adjacency matrix $A = [a_{ij}] \in \mathbb{R}^{n \times n}$, where $a_{ij} = 1$ if missile j is able to communicate directly with missile i , otherwise $a_{ij} = 0$. $\mathcal{N}_i = \{j \in \mathcal{V} : (i, j) \in \xi\}$ is the set of neighboring missiles of the i^{th} missile. In practical engineering, the communication topology is determined through comprehensive considerations of the communication cost and actual demand. In this work, the undirected topology shown in Fig. 2 is adopted, enabling neighboring missiles to share information.

C. Observation

For the multi-missile system, the complete observation information of the entire system is not available to each agent. Thus, the cooperative guidance problem is a partially observable Markov decision process (POMDP) described by

$$O_i \times A_i \rightarrow O'_i, \quad (3)$$

where, O_i and A_i represent the observation and action of the i^{th} missile. O'_i is the observation of the i^{th} missile at next time step.

The full state information of each missile consists of three components: personal features, target features, and error features shown in Table I. P_{mi} and P_T represent the positions of the missile i and the target in two-dimensional coordinates. The target features are estimated or detected through onboard equipment, and the estimation error is assumed to be negligible compared with the required guidance precision. The acquisition of accurate location information requires the support of powerful global positioning systems, here we only need relative error information. e_t^i is the consensus error of time of the missile i :

$$e_t^i = \sum_{j \in \mathcal{N}_i} (t_{go}^i - t_{go}^j), \quad (4)$$

The consensus error of LOS angle of the missile i is defined as:

$$e_a^i = \sum_{j \in \mathcal{N}_i} (e_{\Xi}^i - e_{\Xi}^j), \quad (5)$$

where, $e_{\Xi}^i = \Xi_{mi} - \Xi_{di}$ is the LOS angle error of the missile i , and Ξ_{di} is the desired impact angle of missile i :

$$\Xi_{di} = \Xi_{d1} + \sum_{j=1}^{i-1} \delta_d^j, \quad (6)$$

where, δ_d^i is the desired relative impact angle between two missiles, and Ξ_{d1} is the nominal desired impact angle of the first missile which is determined online. To increase the flexibility and autonomy of the intelligent missile system, the desired Ξ_{di} can be adjusted adaptably instead of being a fixed value.

TABLE I: Full state information of each missile

Features	Symbols
Personal Features	$P_{mi} = (x_i, y_i)$
	α_{mi}
	Ξ_{mi}
	V_{mi}
Target Features	$P_T = (x_T, y_T)$
	V_T
	α_T
Error Features	e_t^i
	e_a^i
	e_{Ξ}^i

D. Fitness evaluation

The reward of each missile at one evaluation step consists of a terminal reward and a flight reward. The objective of the cooperative guidance task is to minimize the error e_t^i , e_a^i , and e_{Ξ}^i . Then, the terminal reward is defined as:

$$r_T^i = (\gamma_a \cdot e^{-\xi_a |e_{\Xi}^i|} + \gamma_t \cdot e^{-\xi_t |e_t^i|}) \cdot \epsilon(k), \quad (7)$$

where, ξ_a , ξ_t , γ_a , γ_t are constant coefficients. $\epsilon(k)$ is the step function defined as

$$\epsilon(k) = \begin{cases} 1, & \text{if } k \text{ is terminal step} \\ 0, & \text{otherwise} \end{cases}. \quad (8)$$

Thus, the terminal reward only reflects the results at the terminal step, and $r_T^i = \gamma_a + \gamma_t$ if and only if $e_{\Xi}^i = 0$ and $e_t^i = 0$. The flight reward is defined as:

$$r_F^i = \beta_a (-1 + e^{-k_a |e_a^i|}) + \beta_t (-1 + e^{-k_t |e_t^i|}), \quad (9)$$

where, k_a , k_t , β_a , and β_t are positive constant coefficients. It can be inferred that $r_F^i \leq 0$ is always true. $r_F^i = 0$ if and only if $e_a^i = 0$ and $e_t^i = 0$. Then, the fitness function of missile i for the cooperative guidance task is defined by

$$F_i = \int_t (r_F^i + r_T^i) dt. \quad (10)$$

Thus, the objective of the cooperative guidance task can be achieved by maximizing the fitness function of each missile.

E. Design of the cooperative guidance law

Based on the requirements of the cooperative guidance task, the guidance law proposed in this paper includes two parts: the tracking control part and the consensus control part. The tracking control part is obtained by proportional navigation guidance(PNG) :

$$u_{pi} = [\beta \dot{\Xi}_{mi} V_{mi}, 0]^T, \quad (11)$$

where, β is the navigation constant. Note that the tracking control part only designs the lateral acceleration.

The consensus control part is modeled by a neural network expressed as

$$\begin{aligned} u_{ei} &= W_{3i}^T \cdot \psi(Z_{2i}), \\ Z_{2i} &= W_{2i}^T \cdot \phi(Z_{1i}), \\ Z_{1i} &= W_{1i}^T \cdot \phi(X_i), \end{aligned} \quad (12)$$

where $W_{3i} \in \mathbb{R}^{q_2 \times 2}$, $W_{2i} \in \mathbb{R}^{q_1 \times q_2}$, and $W_{1i} \in \mathbb{R}^{3 \times q_1}$ denote the weight matrices of the output layer. Z_{1i} and Z_{2i} are the outputs of the first and second hidden layers. q_1 and q_2 are the numbers of neurons in each layer. $\psi(\cdot)$ is the bounded activation function $Tanh(\cdot)$ with $||\psi(\cdot)|| \leq a_{lmax}$, and $\phi(\cdot)$ is the common activation function $Sigmoid(\cdot)$. The input state vector X_i is selected as:

$$X_i = [e_a^i, e_t^i, e_\Xi^i]^T. \quad (13)$$

Thus, the guidance law of the missile i is presented as:

$$u_i = (1 - \eta)u_{pi} + \eta u_{ei}, \quad (14)$$

where η is the guidance gain trading off the tracking control part and the consensus control part.

III. NATURAL CO-EVOLUTIONARY STRATEGY

A. bottleneck of RL

Reinforcement learning is a generic term for a class of value-oriented algorithms. It focuses on solving problems of Markov Decision Process (MDP), which also apply to the guidance problem.

Assume there are n agents in the environment. The joint action set is denoted by $A = A_1 \times A_2 \times \dots \times A_n$, and the system state is denoted by S . At each timestep, each agent takes one step, and with a certain probability, the transition occurs. This can be represented as $S \times A \rightarrow S'$, where S' is the next system state after the transition. The reward is given as $S \times A \times S' \rightarrow \mathbb{R}^n$. In a deterministic environment, the transition probability is 1. For the multi-agent system, the S can be decomposed into individual observations: $S = O_1 + O_2 + \dots + O_n$. A sufficient set of observations must be capable of representing the complete system state. In most cases, the agent does not have access to the complete information of the system. This means they only get partial observations instead of the complete state, making the problem a Partial Observable Markov Decision Process (POMDP).

Two challenges, *Nonstationarity* and *partial observability* [17], impede the research for multi-agent systems. Tons of algorithms have popped up focusing on solving this kind of problem. According to the training process, we can roughly divide them into four types:

- Fully decentralized training, decentralized execution;
- Fully centralized training, decentralized execution;
- Centralized training, centralized execution;
- Value decomposition methods.

The existing works have achieved satisfactory results in coping with problems of less complexity and less requirement for precise control. It has been well-investigated that an ill-distributed value function would seriously stagnate performance. Exploration technologies, such as *Ornstein-Uhlenbeck* noise and stochastic exploration, are used to alleviate this problem.

Applying reinforcement learning (RL) theory, it is possible to handle control tasks with either a single missile [30], [31], or discrete action space [11], [32]. However, for control tasks with multiple agents (missiles), inefficient exploration and non-stationarity can lead to a deterioration of the accuracy of the value function, resulting in either saturated control or coordination loss. Value functions can be advantageous for discrete control, but can be flawed for continuous control tasks with large search spaces. Approaches that constrain the policy space have been discussed in [33], but they heavily rely on prior knowledge and do not scale well to different scenarios. As an alternative, evolutionary strategies have abandoned the use of value functions and have shown the outstanding capability for the aforementioned issues.

B. Natural evolutionary Strategy in multi-agent POMDP

In the evolutionary strategy, individual agent (or its policy) is expressed as a *population*, the group of populations and the environment constitute the *ecosystem*. The objective is to develop the optimal strategy for the group of populations to maximize the *fitnesses* of the ecosystem. For cooperative tasks, the optimal strategy of the ecosystem will be exactly the optimal policy for each population.

$$\arg \max_{u_{tot}^*} F_{tot}(u_{tot}^*) = \begin{pmatrix} \arg \max_{u_1} F_1(u_1) \\ \vdots \\ \arg \max_{u_n} F_n(u_n) \end{pmatrix}, \quad (15)$$

where, u_i which is defined in (14) represents the policy of the i th population and $u_{tot}^* = [u_i^*]_{i=1}^n$ is the joint matrix of individual optimal policy. $F_i(\cdot)$ is its corresponding fitness function and $F_{tot}(\cdot)$ is the joint policy fitness function, more details can be viewed in [34]. However, the inverse is not true:

$$\begin{pmatrix} \arg \max_{u_1} F_1(u_1) \\ \vdots \\ \arg \max_{u_n} F_n(u_n) \end{pmatrix} \neq \arg \max_{u_{tot}^*} F_{tot}(u_{tot}^*). \quad (16)$$

This is because the optimal fitness obtained by one population may be based on the suboptimal fitness obtained by other populations. When the other populations evolve, the previous optima is easy to be broken. To overcome this nonstationary issue, it is best for all populations to evolve simultaneously, that is co-evolution. Each generation updates its parameter at the same time, instead of updating sequentially, mapping in slight variance in fitness values.

C. Optimization in co-evolutionary parameter space

The gradient information is obtained by measuring the contribution of each sample. The parameters of the population are defined as θ , and θ' represents that of the next generation. $p_\psi(\theta'|\theta)$ is the distribution function of θ' under θ , where ψ is the intrinsic parameter. Then the expectation fitness of the next generation is expressed as:

$$\mathbb{E}_{\theta' \sim p_\psi(\theta'|\theta)} F(\theta') = \int_{\theta'} p_\psi(\theta'|\theta) F(\theta') d\theta'. \quad (17)$$

The derivative of Eq. (17) with respect to θ is

$$\nabla_{\theta} \mathbb{E}_{\theta' \sim p_{\psi}(\theta'|\theta)} F(\theta') = \mathbb{E}_{\theta'} \{ \nabla_{\theta'} \log p_{\psi}(\theta'|\theta) F(\theta') \}. \quad (18)$$

If we represent θ' as $\theta + \epsilon$, then we have the similar equation

$$\nabla_{\theta} \mathbb{E}_{\epsilon \sim p_{\psi}(\epsilon)} F(\theta + \epsilon) = \mathbb{E}_{\epsilon} \{ \nabla_{\epsilon} \log p_{\psi}(\epsilon) F(\theta + \epsilon) \}. \quad (19)$$

In an ecosystem with multiple populations, populations will interact and affect the evolutionary process. Thus, the fitness function of the i th population is represented by $F_i(\zeta_i)$, where $\zeta_i = \{\theta_i, \theta_j : j \in \mathcal{N}_i\}$ represents the parameter set of the i th population and its neighboring populations. The expected joint fitness of the next generation is expressed as:

$$\mathbb{E}\{F_i(\zeta'_i)\} = \int_{\zeta'_i} p_{\psi}(\zeta'_i|\zeta_i) F_i(\zeta'_i) d\zeta'_i. \quad (20)$$

where, $p(\zeta'_i|\zeta_i)$ is the joint probability distribution of the next generation over ζ_i . Assume that θ'_i and θ'_j are sampled independently, we have $p(\zeta'_i|\zeta_i) = p(\theta'_i) \prod_{j \in \mathcal{N}_i} p(\theta'_j)$.

The gradient of the joint fitness with respect to θ_i is expressed as

$$\begin{aligned} \nabla_{\theta_i} \mathbb{E}\{F_i(\zeta'_i)\} &= \nabla_{\theta_i} \int_{\zeta'_i} p(\zeta'_i) F_i(\zeta'_i) d\zeta'_i, \\ &= \int_{\theta'_i} \int_{\theta'_j} \cdots \nabla_{\theta_i} p(\theta'_i) F_i(\zeta'_i) \prod_{j \in \mathcal{N}_i} p(\theta'_j) d\theta'_j \prod_{j \in \mathcal{N}_i} d\theta'_j \\ &= \int_{\theta'_i} \cdots \int_{\theta'_i} [\nabla_{\theta_i} \log p(\theta'_i) F_i(\zeta'_i)] p(\theta'_i) \prod_{j \in \mathcal{N}_i} p(\theta'_j) d\theta'_j \prod_{j \in \mathcal{N}_i} d\theta'_j \\ &= \mathbb{E}_{(\zeta'_i)} \{ \nabla_{\theta_i} \log p(\theta'_i) F_i(\zeta'_i) \}. \end{aligned} \quad (21)$$

Note that it has the same format as the version of a single population, it seems to be fine if we just keep the original equation. The influence of θ'_c is counteracted through the calculation of its expectation. However, it is known that the expectation of the joint distribution is approximated through sampling with a limited size. Although individuals are sampled without bias (unbiased estimation), there exists an intrinsic bias for inadequate sampling, and the bias will grow linearly with an increment of distribution dimensionality. So it can be a serious issue when taking the expectation of all neighboring parameters, and the sample size stays relatively small.

It is not necessary to take account of all parameters since only the expectation of θ_i is actually needed. To alleviate the incremental bias, we propose to approximate only the expectation of the parameter of the current population θ_i and ignore its neighbor parameters, which is

$$\mathbb{E}_{\theta'_i} F_i(\theta'_i) = \int_{\theta'_i} F_i(\theta'_i) p(\theta'_i) d\theta'_i. \quad (22)$$

Though $p(\theta'_i)$ is available for independent distribution, it is infeasible to obtain $F_i(\theta'_i)$, since all agents are sampled and

evaluated together. However, the expectation of individual fitness can be approximated by the multiplication between the original fitness $F_i(\zeta'_i)$ and its confidence. The rectified expectation is expressed as

$$\mathbb{E}_{\theta'_i} F_i(\theta'_i) = \int_{\theta'_i} F_i(\zeta'_i) p(\theta'_i) \prod_{c \in \mathcal{N}_i} p(\theta'_c) d\theta'_i, \quad (23)$$

where, $\prod_{c \in \mathcal{N}_i} p(\theta'_c)$ is the confidence, and θ'_c represents the samples that appear along with θ'_i . In this way, the bias of estimating the expectation of the neighboring distributions is addressed. The gradient after modification is

$$\begin{aligned} \nabla_{\theta_i} \mathbb{E}_{\theta'_i} F_i(\theta'_i) &= \int_{\theta'_i} [\nabla_{\theta_i} \log p(\theta'_i) F_i(\zeta'_i)] p(\theta'_i) \prod_{c \in \mathcal{N}_i} p(\theta'_c) d\theta'_i \\ &= \mathbb{E}_{\theta'_i} \left\{ \nabla_{\theta_i} \log p(\theta'_i) F_i(\zeta'_i) \prod_{c \in \mathcal{N}_i} p(\theta'_c) \right\}. \end{aligned} \quad (24)$$

Remark 1. We refer to strategies that use (24) as the updating policy as the natural co-evolutionary strategy (NCES). The natural co-evolutionary strategy is a strategy for evolutionary algorithms that updates the weights of the population in multi-agent systems. Compared to other existing evolutionary strategy algorithms, NCES alleviates the incremental bias caused by neighboring parameters and achieves better performance in cooperative continuous control tasks.

The core idea is that although the individual fitness $F_i(\theta'_i)$ does not exist, the expectation of the individual fitness does, and is invariant to the parameter distributions of its neighboring agents, so the expectation of the individual agent's fitness should be calculated instead of including the expectation of neighboring agents. Let's denote the expectation of the objective function over (ζ'_i) , which is $\nabla_{\theta_i} \log p(\theta'_i) F_i(\zeta'_i)$, by $\phi(\zeta'_i)$ and the expectation of the objective function over θ'_i by $\phi(\theta'_i)$, such that

$$\phi(\theta'_i) = \int_{\theta'_c} \phi(\zeta'_i) p(\theta'_c) d\theta'_c. \quad (25)$$

To visualize the sampling estimation process, we use a variant of *eggholder* as the objective function for demonstration, which is defined in (26), since the real objective is too expensive to obtain.

Assume there exists one neighboring population θ_c for θ_i with the size of 400, the sampled individuals are shown in Fig. 3, following a bivariate normal distribution, and the parameter spaces are confined to $\theta'_i \in [-2, 2]$, $\theta'_c \in [-2, 2]$. Since θ'_i and θ'_c are sampled independently, the individuals can be considered to be sampled from $p(\theta'_i)$ only, which is represented by the sample points in Fig. 4. In this objective graph with single-dimensional parameter space, the real objective curve expressed in solid line is obtained by (25). In order to standardize the scope, all the sampled data including the real objective values are uniformly scaled to the range $[0, 1]$, and such standardization does not affect the directionality of the estimated gradient.

The original objective value $\phi(\theta'_i)$ for each sample varies as the corresponding $p(\theta'_c)$ changes, which introduces additional estimation bias. As shown by the blue dots in Fig. 4, the

TABLE II: A variant of the eggholder function

$$\phi(\zeta'_i) = \frac{(30\theta'_c + 47) \sin \sqrt{|30\theta'_c + 15\theta'_i + 47|} - 30\theta'_i \sin \sqrt{|30\theta'_i - (30\theta'_c + 47)|}}{200} - 0.2 \quad (26)$$

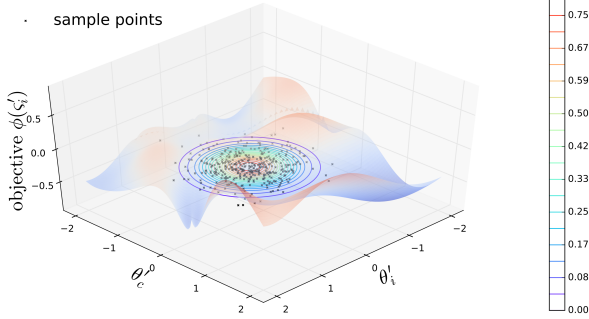


Fig. 3: Estimate gradient through sampling

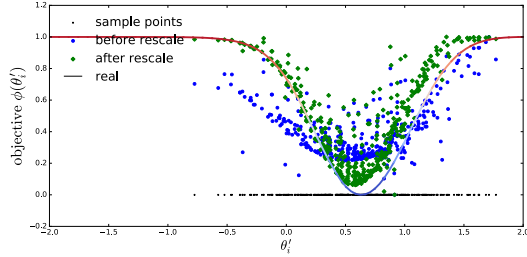


Fig. 4: Estimate gradient through sampling from single distribution

distribution of the objective values before rescale is significantly different from the distribution of the true objective values. From Fig. 3, it can be seen that as sample points deviate from the distribution center, their probability of being sampled also decreases, which means that the accuracy or confidence of the fitness of each sample $\phi(\zeta'_i)$ decreases with the decrease of $p(\theta'_c)$. If the original objective is rescaled by its confidence $p(\theta'_c)$, which is the probability of the appearance of the ζ'_i given the existence of θ'_i , the reconstructed objective values represented by the green square dot in Fig. 4 is closer to the real $\phi(\theta'_i)$, which obviously reduce the estimation bias.

The above proof indicates that in the case of limited population size and a large number of neighboring populations, applying the rescaled gradient will keep the approximation bias to the level of a single population, resulting in a more accurate estimation of gradient information, empirical results also supported this conclusion. However, when the population size is large enough (e.g., thousands), this approach may not result in additional accuracy improvements.

The modified expression is also desirable for parallel computing, as only the perturbation of the neighboring populations is needed, which can be easily obtained through communication among processes, and the probabilities can be calculated in a distributed approach.

D. Elitist adaptation Techniques

The performance of NES is sensitive to hyper-parameters, and the learning rate is usually the most critical hyper-parameter of NES. Thus, an elitist adaptation method for the learning rate is applied in this paper. First, a list of learning rates is linearly selected in the neighborhood of the original learning rate η_α as:

$$\eta_{cad} = \{clip((1 + 0.1k)\eta_\alpha, \eta_{\alpha min}, \eta_{\alpha max}) : k \in \mathbb{Z}, -l/2 \leq k \leq l/2\}, \quad (27)$$

where, $\eta_{cad} \in \mathbb{R}^{m+1}$. The $\eta_{\alpha min}$ and $\eta_{\alpha max}$ are the minimum and maximum value of η_α . l is the size of perturbations which is clipped by $clip(\cdot)$. To evaluate the quality of the candidate learning rates, the evaluation function $G_i(\cdot)$ is defined:

$$G_i(\eta_{cad}) = \begin{pmatrix} F_i(\theta_i + \eta_{cad}^{-l/2} g_{\theta_i}) - F_i(\theta_i + \eta_\alpha g_{\theta_i}) \\ \vdots \\ F_i(\theta_i + \eta_{cad}^{l/2} g_{\theta_i}) - F_i(\theta_i + \eta_\alpha g_{\theta_i}) \end{pmatrix}, \quad (28)$$

where η_{cad}^k is the k th sampled learning rate of the candidate list. The gradient g_{θ_i} is kept after evaluation. Therefore, by comparing the candidate learning rates with the original one, the next update can be better than the previous one. Considering *peer pressure*, each missile is assigned the same learning rate. The learning rate of the next generation is obtained by

$$\eta'_\alpha = \arg \max_{\eta_{cad}^k} \left(\sum_{i=1}^n G_i(\eta_{cad}) \right). \quad (29)$$

A similar approach is employed to obtain the optimal Ξ_{d1}^* during the training process.

$$\Xi_{d1}^* = \arg \max_{\Xi_{d1}^k} \begin{pmatrix} H(\Xi_{d1}^1) \\ \vdots \\ H(\Xi_{d1}^h) \end{pmatrix} \quad (30)$$

where, Ξ_{d1}^k is uniformly sampled from the region $[-\pi, \pi]$. $H(\cdot)$ is the fitness function of sampled LOS angle that is defined as

$$H(\Xi_{d1}^k) = F_{tot}(\theta_{init})|_{\Xi_{d1}=\Xi_{d1}^k}, \quad (31)$$

where $\theta_{init} = [\theta_i^{init}]$ is the joint initial individual parameters. In this way, the desired impact angles are established automatically.

A rank-based fitness shaping method that is in the same spirit as the one proposed in [20] is employed in shaping the raw fitness. Conventionally, we still let $F_i(\cdot)$ denote the fitness function after shaping. Another technique called mirrored sampling [18] is also applied for sampling parameter perturbations.

IV. HYBRID CO-EVOLUTIONARY COOPERATIVE GUIDANCE ALGORITHM

To achieve coordinated attack, the natural co-evolutionary strategy is applied to optimize the parameter matrices $\theta_i = [W_{3i}, W_{2i}, W_{1i}]$ of the neural network controller.

The univariate Gaussian distribution with zero means and standard deviation σ is used to sample perturbations. According to (24), it can be obtained that:

$$g_{\theta_i} = \mathbb{E}_{\epsilon_i \sim N(0, \sigma^2)} \left\{ \nabla_{\theta_i} \log p(\theta'_i) F_i(\zeta'_i) \prod_{c \in \mathcal{N}_i} p(\theta'_c) \right\} \quad (32)$$

$$= \frac{1}{m\sigma^2} \sum_{i=1}^m F_i(\zeta'_i) \epsilon_i \prod_{c \in \mathcal{N}_i} p(\epsilon_c).$$

The complete implementation algorithm of the proposed guidance law is shown in Algorithm 1. The conceptual diagram in Fig. 5 figuratively revealed the parallel simulation process. A master-slave (or fully-distributed) model [35] [36] is used for large-scale parallel computation. In this case, each population is evaluated in a separate process and the results of the ecosystem are aggregated to calculate the rescaled gradient (32) and sent to produce guided generations. The sampled generations are then distributed to each parallel process, and the gradient is recalculated and updated.

Algorithm 1 Hybrid Cooperative Co-Evolutionary Guidance Law (HCCGL)

Require: $\eta_\alpha, \eta, \sigma, \theta_{init} = [\theta_i^{init}]$, agent number n .
 Sample $[\Xi_{d1}^k] \in \mathbb{R}^h \sim U(-\pi, \pi)$, obtain Ξ_{d1}^* using (30)
repeat
 for $k = 1 \dots m$ **do**
 Sample group of individuals:
 $\epsilon^k = \{\epsilon_i^k \sim N(0, \sigma^2 I) : i \in \{1, \dots, n\}\}$,
 $\zeta^k = \{\zeta_i^k = \{\theta_i^k, \theta_j^k : j \in \mathcal{N}_i\} : i \in \{1, \dots, n\}\}$
 evaluate fitness $F_i(\zeta_i^k)$, for $i \in \{1, \dots, n\}$
 end for
 for each agent $i = 1 \dots n$ **do**
 calculate natural gradient:
 $g_{\theta_i} \leftarrow \frac{1}{m\sigma^2} \sum_{k=1}^m F_i(\zeta_i^k) \epsilon_i^k \prod_{c \in \mathcal{N}_i} p(\epsilon_c^k)$
 $\theta_i \leftarrow \theta_i + \eta_\alpha \cdot g_{\theta_i}$
 end for
 if time for adaptation **then**
 sample η_{cadi} using (27)
 $\eta_\alpha \leftarrow \arg \max_{\eta_{cad}} (\sum_{i=1}^n G_i(\eta_{cad}))$
 end if
until stopping criterion is met

Theorem 1. Under the control policy (14) and the update strategy shown in Algorithm 1, by selecting an appropriate learning rate, sampling variance, and population size, the obtained control policy will converge to a small neighborhood of the optimal control policy when $l \rightarrow \infty$.

Proof. The approximation error of the control policy at the l th iteration is defined by

$$E_{ui}^{[l]} = u_i^{[l]} - u_i^*, \quad (33)$$

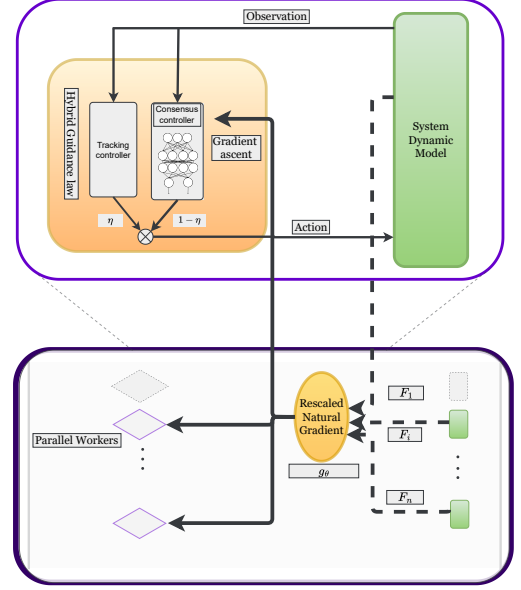


Fig. 5: Conceptual framework of our proposed HCCGL, the upper box connected by dotted lines is a detailed expansion of the evaluation and evolutionary processes in the lower box.

where u_i^* is the optimal control policy of the i th agent, and we have

$$u_i^{[l]} = (1 - \eta)u_{pi} + \eta u_{ei}^{[l]}. \quad (34)$$

The control policy of the neural network is represented by its parameter set. Since a neural network with a single hidden layer can approximate a multivariate continuous function with arbitrary precision [37], which implies that a single hidden layer perceptron with sufficient units is equivalent to the neural network with three hidden layers used in this work. In this way, the neural network controller can be represented by

$$u_{ei}^{[l]} = W_i^{[l]T}, \quad (35)$$

which is a column matrix and the activation function parameters are regarded as constants. By combining (35) and (34) and substituting it into (33), we have

$$E_{ui}^{[l]} = (1 - \eta)u_{pi} + \eta W_i^{[l]T} - u_i^*, \quad (36)$$

$$E_{ui}^{[l-1]} = (1 - \eta)u_{pi} + \eta W_i^{[l-1]T} - u_i^*. \quad (37)$$

Further combining (36) and (37), the term of fixed controllers are eliminated, and we have

$$\begin{aligned} E_{ui}^{[l]} - E_{ui}^{[l-1]} &= \eta(W_i^{[l]T} - W_i^{[l-1]T}) \\ &= \eta \eta_\alpha g_{\theta_i}^{[l-1]}, \end{aligned} \quad (38)$$

where η_α is the learning rate and η the guidance gain, with $\eta_\alpha, \eta > 0$. The focus of this equation, $g_{\theta_i}^{[l-1]}$, is the policy

update gradient at the l -th iteration, which follows by (32). Expanding this equation, we have

$$\begin{aligned}
 E_{ui}^{[l]} - E_{ui}^{[l-1]} &= \eta \eta_\alpha g_{\theta_i}^{[l-1]} \\
 &= \eta \eta_\alpha \frac{1}{m\sigma^2} \sum_{k=1}^m F_i(E_{ui}^{[k]}) \epsilon_i^k \prod_{c \in \mathcal{N}_i} p(\epsilon_c^k) \\
 &= \eta \eta_\alpha \frac{1}{m\sigma^2} \sum_{k=1}^m F_i(E_{ui}^{[k]}) * (E_{ui}^{[k]} - E_{ui}^{[l-1]}) \\
 &\quad \prod_{c \in \mathcal{N}_i} p(\epsilon_c^k).
 \end{aligned} \tag{39}$$

Note that in (39), $F_i(E_{ui}^{[k]}) : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is the transformed fitness function for the evaluation of the policy error, with p as the number of parameters. Thus, it is different from the fitness function discussed in the previous sections, which evaluates the policy directly. It is assumed that $F_i(E_{ui}^{[k]})$ is fully differentiable to the policy controller, and

$$\frac{\partial F_i(|E_{ui}^{[k]}|)}{\partial |E_{ui}^{[k]}|} < 0, \tag{40}$$

considering that $|E_{ui}^{[k]}|$ represents the quality of the policy globally.

In an effort to linearize the fitness evaluation function, Taylor's formula is utilized to expand the equation at $|E_{ui}^{[l-1]}|$ and the higher order terms are ignored, and we obtain

$$F_i(|E_{ui}^{[k]}|) = G_i^{[l-1]}[|E_{ui}^{[k]}| - |E_{ui}^{[l-1]}|] + F_i(|E_{ui}^{[l-1]}|), \tag{41}$$

where $G_i^{[l-1]} \in \mathbb{R}^{p \times p}$ is a diagonal Jacobian matrix defined by

$$G_i^{[l-1]} = \left. \frac{\partial F_i(|E_{ui}^{[k]}|)}{\partial |E_{ui}^{[k]}|} \right|_{|E_{ui}^{[l-1]}|}, \tag{42}$$

with negative entries and p as the number of total parameters. Since $|E_{ui}^{[k]}|$ is located within a tiny vicinity of $|E_{ui}^{[l-1]}|$, (41) is of considerable accuracy.

Then, by taking the absolute value of the approximation error and substituting (41) into (39), and considering $\Delta E_i^k = |E_{ui}^{[k]}| - |E_{ui}^{[l-1]}|$ we obtain

$$\begin{aligned}
 \Delta |E_{ui}^{[l]}| &= |E_{ui}^{[l]}| - |E_{ui}^{[l-1]}| \\
 &= \eta \eta_\alpha \frac{1}{m\sigma^2} \sum_{k=1}^m [G_i^{[l-1]} \Delta E_i^k + F_i(|E_{ui}^{[l-1]}|)] * \\
 &\quad \Delta E_i^k * \prod_{c \in \mathcal{N}_i} p(\epsilon_c^k) \\
 &= \eta \eta_\alpha \frac{1}{m\sigma^2} G_i^{[l-1]} \sum_{k=1}^m [\Delta E_i^k * \Delta E_i^k + \\
 &\quad F_i(|E_{ui}^{[l-1]}|) * [|E_{ui}^{[k]}| - |E_{ui}^{[l-1]}|]] * \prod_{c \in \mathcal{N}_i} p(\epsilon_c^k) \\
 &= \eta \eta_\alpha \frac{1}{m\sigma^2} [G_i^{[l-1]} \sum_{k=1}^m \Delta E_i^k * \Delta E_i^k * \prod_{c \in \mathcal{N}_i} p(\epsilon_c^k) \\
 &\quad + F_i(|E_{ui}^{[l-1]}|) * \sum_{k=1}^m \Delta E_i^k * \prod_{c \in \mathcal{N}_i} p(\epsilon_c^k)].
 \end{aligned} \tag{43}$$

For brevity, we define A_i^k , P_i^k , and B_i^k by

$$\begin{aligned}
 A_i^k &= \Delta E_i^k * \Delta E_i^k; \\
 P_i^k &= \prod_{c \in \mathcal{N}_i} p(\epsilon_c^k); \\
 B_i^k &= \sum_{k=1}^m |E_{ui}^{[k]}| * P_i^k - \sum_{k=1}^m |E_{ui}^{[l-1]}| * P_i^k,
 \end{aligned} \tag{44}$$

such that

$$\begin{aligned}
 \Delta |E_{ui}^{[l]}| &= \eta \eta_\alpha \frac{1}{m\sigma^2} [G_i^{[l-1]} \sum_{k=1}^m A_i^k * P_i^k + \\
 &\quad F_i(|E_{ui}^{[l-1]}|) * B_i^k].
 \end{aligned} \tag{45}$$

Since $|E_{ui}^{[k]}|$ is sampled from an unbiased normal distribution which is centered at $|E_{ui}^{[l-1]}|$, as shown in the analysis of Section III-C, we have

$$B_i^k \rightarrow 0, \quad \text{as } m \rightarrow \infty. \tag{46}$$

Also, from the matrix Hadamard product we have

$$A_i^k > 0, \tag{47}$$

and

$$P_i^k > 0. \tag{48}$$

$G_i^{[l-1]}$ is negative definite. Given sufficient large m , it is evident that

$$\Delta |E_{ui}^{[l]}| < 0, \quad l = 1, 2, \dots \tag{49}$$

Therefore, by adjusting the learning rate η_α attentively, the approximation error can be decreased to a considerably small range δ_e , such that

$$\begin{aligned}
 \lim_{l \rightarrow \infty} |E_{ui}^{[l]}| &= \delta_e, \\
 \delta_e &\rightarrow 0, \quad \text{as } m \rightarrow \infty.
 \end{aligned} \tag{50}$$

Thus, the control policy u_i converges to a small neighborhood of the optimal control policy u_i^* , resulting in a stabilizing control system. \square

V. SIMULATIONS AND ANALYSIS

To verify the validity of the proposed method, a variety of simulations based on the cooperative guidance framework are designed. Both cases with stationary target and maneuvering target are simulated. Further, comparison experiments are performed to fully demonstrate the superiority of the proposed guidance method.

A. Parameter setup

The acceleration constraint and velocity constraint of the missiles are listed in Table III. The hyper-parameters of the algorithm are listed in Table IV.

Now that frameskip has been extensively employed in continuous control problems [19]. In this work, this parameter of frameskip is set to 12 for case 1 and case 2, and 40 for case 3. Appropriate adjustment of this parameter will facilitate the training process without affecting the final results.

TABLE III: Constraints of the missiles

Parameter	Value
maximum lateral overload (g), a_{lmax}	50
maximum trust overload (g), a_{vmax}	5
Upper bound of velocity (m/s), V_{max}	900
Lower bound of velocity (m/s), V_{min}	350

TABLE IV: Hyper-parameters of the cooperative guidance algorithm

Parameter	Value
simulation step (ms), τ	5
guidance gain, η	0.3
Initial learning rate, η_α	0.015
standard deviation for sampling population, σ	0.2
size of learning rate adaptation, l	20
size of population, m	140
adaptation cycle, ρ	50
navigation constant, β	4
k_a	1
k_t	0.2
ξ_a	10
ξ_t	1
λ_a	4000
λ_t	2000
β_a	10
β_t	2

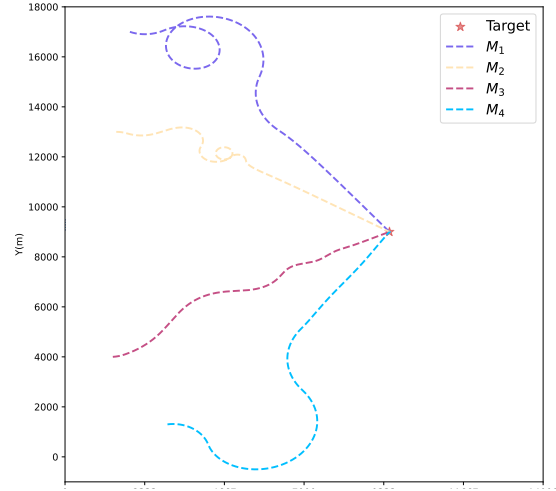
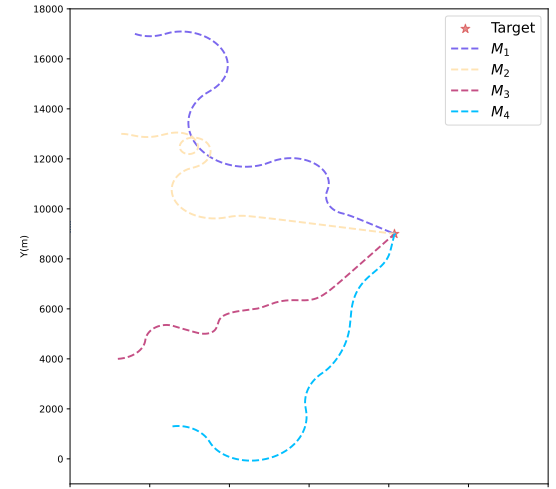
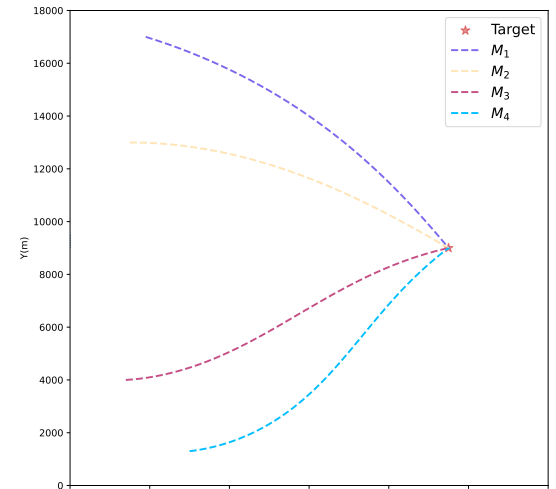
B. Case 1: Comparison Experiments

In this section, the proposed guidance law is compared with the time and space cooperative guidance law (TASCGL) proposed in [38], which considers the space and time cooperative guidance under the distributed communication topology. However, different from the method proposed in this work, the compared method is susceptible and brittle to the initial conditions. Therefore, in order to verify the generalization ability of the control methods, a uniform initial condition was adopted in the comparison simulation, which differs slightly from the initial condition in the comparison method. The initial conditions as shown in Table V. Four missiles are engaged in the cooperative scenario with different desired relative impact angles δ_d^i as 20° , 60° , and 30° , for each $i = 1, 2, 3$, respectively. The target is located at (9500, 9000)m.

Although the reference method is primarily designed for directed topology, it can be well extended to an undirected topology condition, thus in order to conduct effective comparison experiments, we additionally implemented a comparison experiment under an undirected topology the same as the one used in the proposed method. We use TASCGL^a and TASCGL^b to denote the comparative experiments performed under directed and undirected communication topologies, respectively.

TABLE V: Initial conditions of case 1.

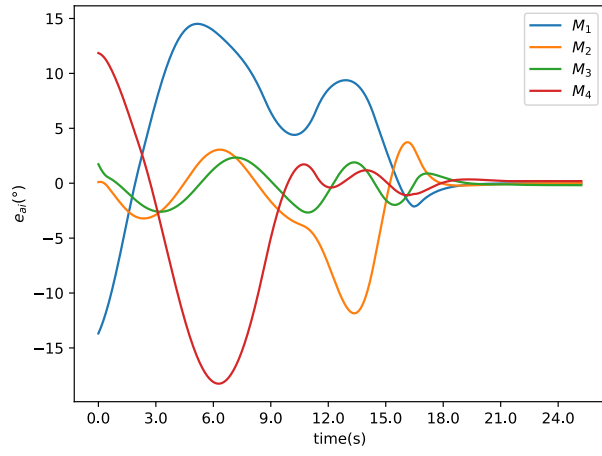
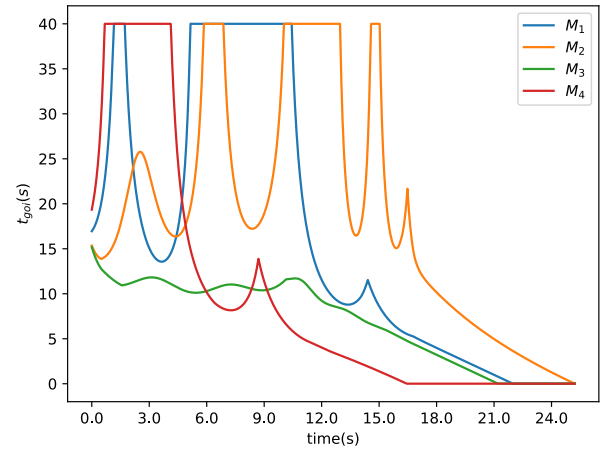
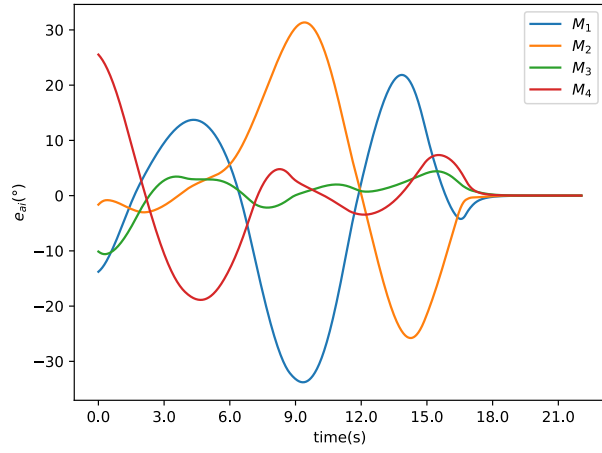
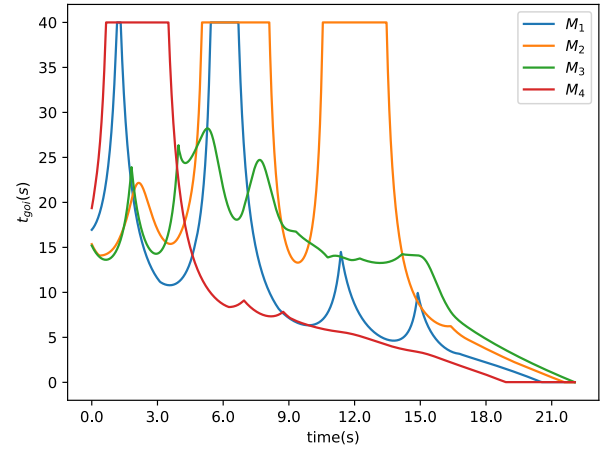
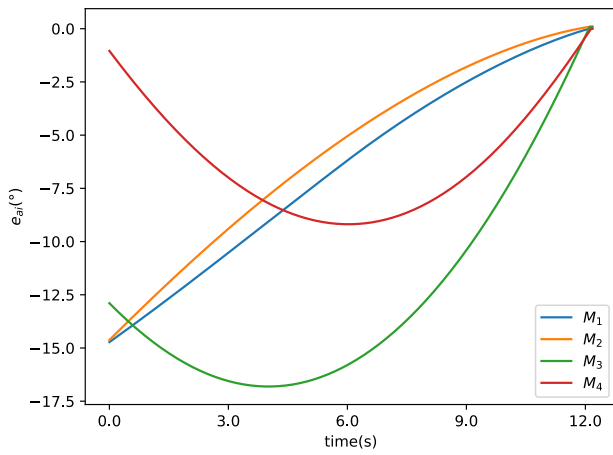
Missile	Position (m)	Flight-path Angle ($^\circ$)	Velocity (m/s)
M_1	(1900, 17000)	-25	700
M_2	(1500, 13000)	0	650
M_3	(1400, 4000)	5	700
M_4	(3000, 1300)	10	680

(a) TASCGL^a(b) TASCGL^b

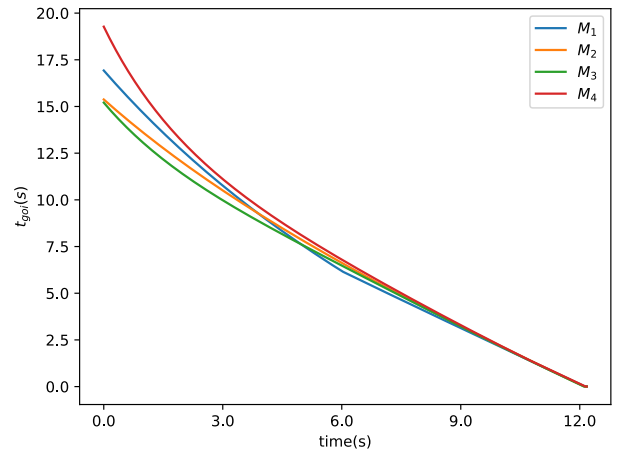
(c) HCCGL

Fig. 6: Trajectories of the two methods

Fig. 6 shows the trajectories of the two guidance laws. As depicted in the figure, the trajectory of TASCGL is twisted at the initial stage, as the missiles try to consensus their LOS angles and velocities. In comparison, the trajectory of the

(a) TASCGL^a(a) TASCGL^a(b) TASCGL^b(b) TASCGL^b

(c) HCCGL



(c) HCCGL

Fig. 7: Consensus angle error profiles of the two methods

Fig. 8: Time-to-go profiles of the two methods

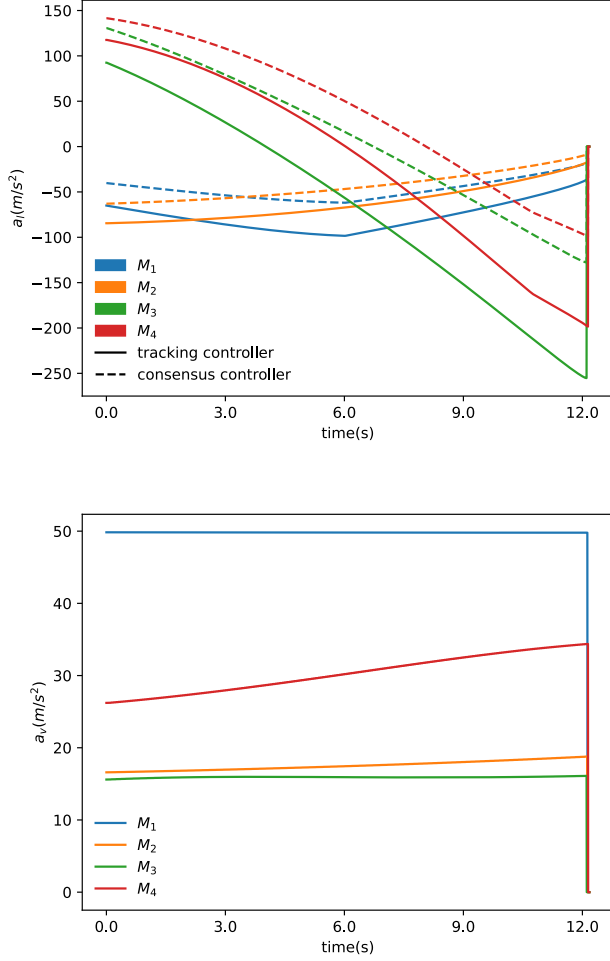


Fig. 9: Decomposition of acceleration commands in Case 1

proposed HCCGL exhibited better damping performance with no oscillations.

It can be seen from Table VI that the Zero-Effort Miss (ZEM) and the consensus angle error for both guidance laws have achieved competitive final accuracy. The consensus time error of TASCGL was up to 5 seconds under both directed and undirected topologies, whereas the proposed method achieved an error of less than 0.1 seconds. Further analysis of the velocity curve shows that in the case of TASCGL, the velocities are prohibited from reaching their ideal values due to the velocity boundary, which is not considered in its design, thus leading to desynchronization in impact time. The profiles of the two methods are shown in Fig. 7 and Fig. 8, it can be observed that the flight time of all missiles under HCCGL trends to be identical. For HCCGL, the decomposition of acceleration commands is shown in Fig. 9. The left figure shows the decomposition of lateral accelerations, in which the solid line represents the command from the tracking controller while the dashed line represents the command from the consensus controller before weighing. Since the tracking part is derived from proportional navigation, the vertical acceleration shown on the right one is completely derived from the consensus

controller. The two parts of accelerations have similar trends but do not coincide, demonstrating the effectiveness of the consensus controller, which is trained with the improved co-evolutionary strategy.

The result reveals that the proposed guidance law outperforms the compared method with higher precision in consensus performance and smoother trajectories. Moreover, as the traditional guidance law is usually constrained to boundary conditions and missile's superb maneuverability, the proposed guidance law is more resilient to limited conditions and more intelligent to be aware of the time-varying states of missiles of collaboration.

TABLE VI: Comparison results of two guidance laws in case 1.

Algorithm	Index	M_1	M_2	M_3	M_4
TASCGL ^a	$e_t^i(s)$	5.54	3.23	-4.02	-4.75
	$e_a^i(^{\circ})$	-9.83E-3	-1.58E-3	-1.80E-1	1.91E-1
	ZEM(m)	2.24E-7	5.06E-7	-7.00E-2	3.00E-4
TASCGL ^b	$e_t^i(s)$	6.15E-1	5.5E-1	3.655	-4.82
	$e_a^i(^{\circ})$	-6.11E-3	-1.28E-3	-3.38E-4	7.73E-3
	ZEM(m)	1.19E-5	5.11E-8	4.81E-7	5.53E-5
HCCGL	$e_t^i(s)$	-1.00E-2	1.00E-2	-5.00E-2	5.00E-2
	$e_a^i(^{\circ})$	1.79E-2	9.67E-2	9.20E-2	4.69E-3
	ZEM(m)	4.19E-5	7.66E-6	3.11E-4	4.09E-5

C. Case 2: Non-stationary target

In this part, an engagement scenario with a non-stationary target is designed and simulated to verify the effectiveness of the proposed method against unknown dynamic target. The target is maneuvering with lateral acceleration $a_t = 5g \sin(\frac{\pi}{7}t)$ with its velocity fixed at $V_t = 130m/s$, and its initial flight-path angle $\alpha_T = 162^{\circ}$. Other initial conditions are the same in case 1. Simulation trajectory and the result can be seen in Fig. 10 and Table VII.

From Table VII we can see that the consensus angle error is within one degree, which is sufficient for the accuracy requirement, and salvo attack is achieved with negligible consensus time error. The result demonstrates the effectiveness of the proposed guidance method in intercepting the dynamic target. As far as the author knows, it is the first time achieving cooperative guidance against non-stationary target with intelligent control, which shows its extraordinary robustness against disturbance from non-stationary objectives.

TABLE VII: Result for Case 2.

Index	M_1	M_2	M_3	M_4
$e_a^i(^{\circ})$	-3.43E-1	-6.53E-2	-1.23E-1	-1.49E-1
$e_t^i(s)$	6.00E-2	1.85E-2	2.10E-1	-4.55E-1
ZEM(m)	1.83E-3	6.38E-3	2.49E-2	9.28E-1

D. Case 3: Monte-Carlo simulation

Monte-Carlo simulation has been extensively employed to examine the robustness of an algorithm under varying initial conditions, thus it is applied in this section. In the

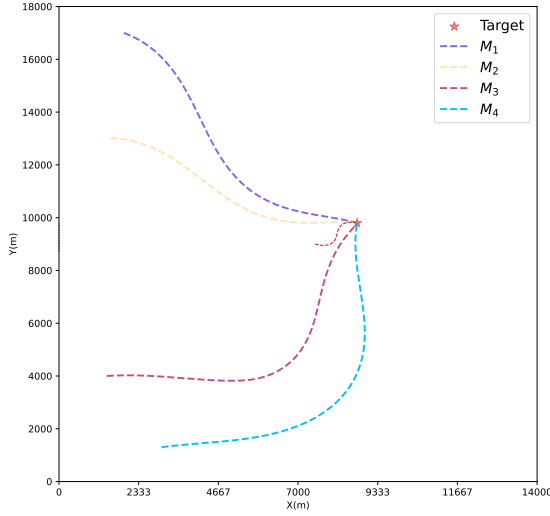


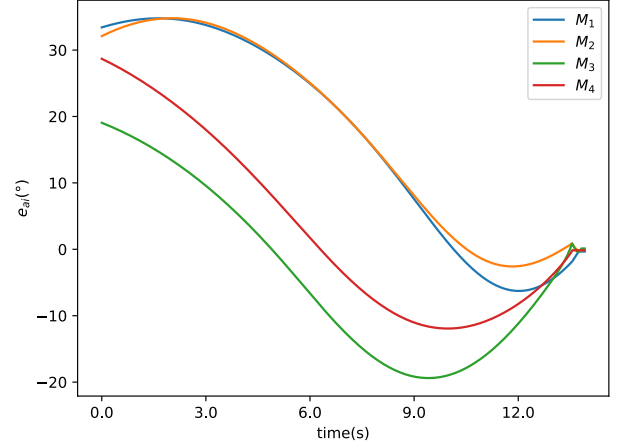
Fig. 10: Trajectory in case 2

existing literature, the target is usually regarded as stationary as interception of a stationary target is more exclusive of unpredictable disturbance. In this case, five missiles are engaged, and each missile's position is randomly sampled from a uniform distribution, which is denoted by $U(\cdot, \cdot)$. Specifically, for the i^{th} missile, the x-coordinate of its position is $U(2000, 2600)$ and the y-coordinate is $U(11000, 13000) - 2000i$, which makes the missiles arranged in an orderly manner. The initial flight-path angles of all missiles are set to 0° , with identical velocities of 600m/s and the same desired relative impact angles of 25° . Additionally, the target's position is (10000m, 9000m).

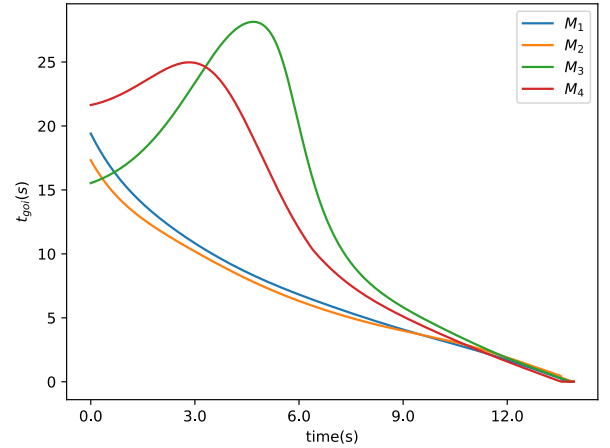
Simulations with randomly sampled conditions are conducted in 200 episodes. The diverse trajectories are depicted in Fig. 12, and the statistical result after taking the absolute value is shown in Table VIII. From the result, we can see that the mean errors of impact angles are within 1° , and the consensus error of impact time holds within 1s most of the time. The result shows that for any initial state with limited error, the proposed scheme can always find the relative optimal solution.

TABLE VIII: Result for Case 3.

Index		M_1	M_2	M_3	M_4	M_5
$e_a^i(^{\circ})$	Mean	4.50E-1	8.20E-1	7.10E-1	2.20E-1	9.30E-1
	Max	1.85E-0	3.37E-0	1.96E-0	6.10E-1	2.37E-0
	Min	4.56E-3	6.40E-3	7.01E-4	4.58E-4	6.46E-3
$e_t^i(s)$	Mean	6.10E-1	5.50E-1	5.30E-1	4.50E-1	5.50E-1
	Max	1.78E-0	1.57E-0	1.63E-0	1.54E-0	1.44E-0
	Min	1.50E-2	1.00E-2	1.78E-15	5.00E-3	1.78E-15
ZEM(m)	Mean	5.85E-3	5.89E-3	3.74E-4	9.05E-4	9.93E-4
	Max	1.03E-2	1.07E-2	7.77E-4	2.53E-3	3.39E-3
	Min	2.18E-3	2.04E-3	6.42E-5	1.68E-5	2.72E-6



(a) Consensus angle error profile



(b) Time-to-goes profile

Fig. 11: Flight data profiles in Case 2

E. Optimization process analysis

Fig. 13 shows the learning curves in the three cases. The mean fitness in case 1 keeps moving upper and merges together at the final phase. From the curve of case 2, we can see that two of the missiles get ahead about 1000 scores, but finally back to meet with the other missiles. A similar phenomenon also appears in case 3. It can be inferred that the policies asymptotically evolved to the equilibrium state, and one reason is that the rescaled gradient prohibited the ever-increasing gap between individual groups, which is crucial for mutual improvement. If one group gets ahead too much, then the other groups may never chase up due to the interrelationship, which is to say that the improvement of the poorer-performed group is prohibited when more significant drops in the better-performed ones will occur. Fig. 14 presents the adaptation profiles of learning rates applying the aforementioned technique. For case 1 and case 2, the learning rates start from high values and gradually converge to the minimal value, which corresponds with the quality of estimated gradients. However, due to the

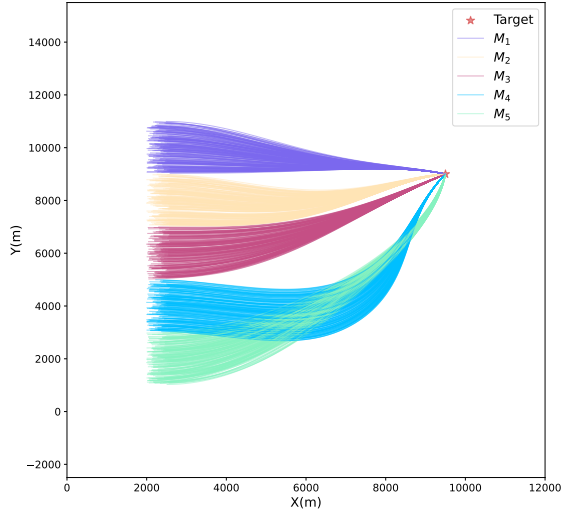


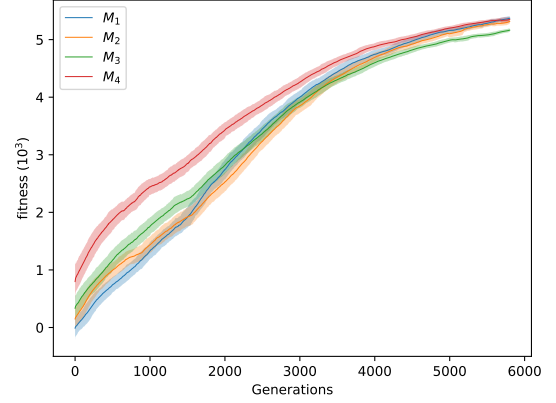
Fig. 12: Diverse trajectories of the Monte-Carlo simulation

random initial conditions in case 3, the learning rates will not settle easily. The extensive empirical result shows that without the learning rate adaptation, the fitness profiles will jitter in the end instead of converging to satisfactory ranges (regardless of the types of optimizer). Note that it is pretty common when training neural networks and may presumably have been caused by overfitting, according to related research in the field. Employing the simple adaptation technique contributes to distressing this deficiency.

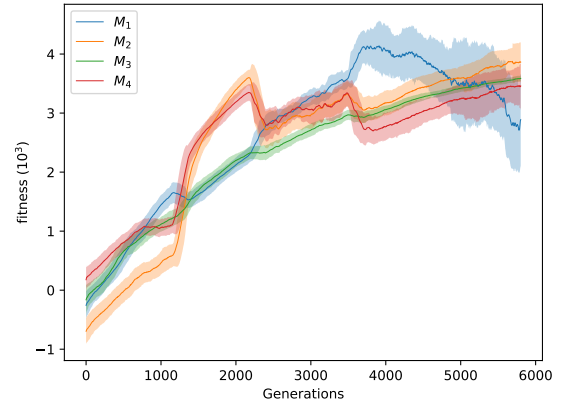
VI. CONCLUSIONS

In this paper, an improved co-evolutionary strategy NCES has been developed to solve the non-stationarity issue in multi-agent dynamic environments. The hybrid co-evolutionary cooperative guidance law (HCCGL) has been proposed to integrate with the improved strategy, and the neural network has been used to construct the consensus controller. To fully demonstrate its effectiveness in synchronizing impact time and angles, three experiments under different conditions have been carried out. Experiment on maneuvering target has been proven effective with satisfactory precision. The proposed method is shown to be robust and can be well scaled to solve the cooperative guidance problem for the multi-agent system, which is the first time an intelligent cooperative guidance law is applied to intercept a non-stationary target with time and angle constraints in the existing studies.

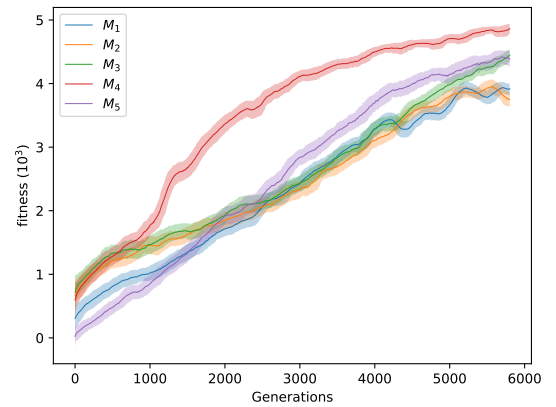
The proposed algorithm combines traditional control theories with intelligent algorithms, revealing the enormous potential in this field. It is always meaningful to explore the limits of modern control tasks. Despite the satisfactory results that have been acquired, this work still left space to be improved. Future works may include exploring the effectiveness of incremental guidance gain, or control strategies that tackle actuation failure and system uncertainty.



(a) Case 1

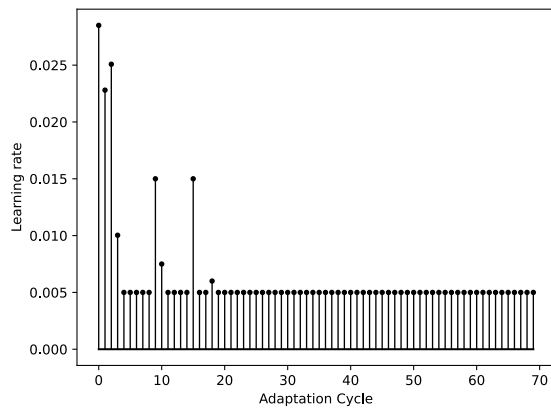


(b) Case 2

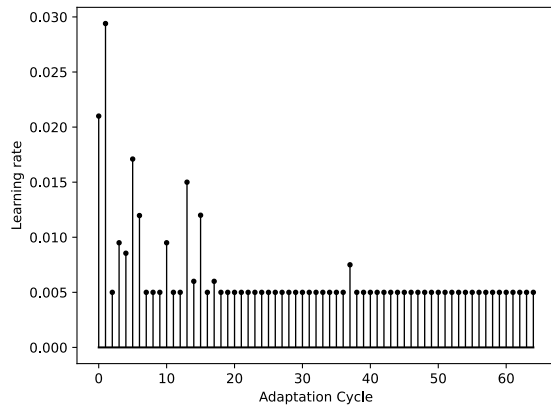


(c) Case 3

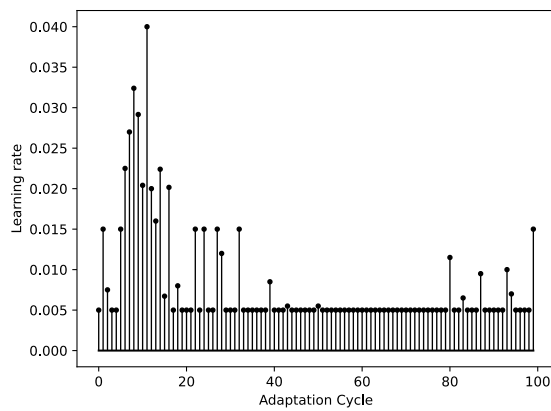
Fig. 13: Mean fitness profiles over iterations of three cases



(a) Case 1



(b) Case 2



(c) Case 3

Fig. 14: Learning rate profiles over evaluation iterations of three cases

REFERENCES

- [1] I.-S. Jeon, J.-I. Lee, and M.-J. Tahk, "Homing guidance law for cooperative attack of multiple missiles," *Journal of guidance, control, and dynamics*, vol. 33, no. 1, pp. 275–280, 2010.
- [2] K. Ma, H. K. Khalil, and Y. Yao, "Guidance law implementation with performance recovery using an extended high-gain observer," *Aerospace Science and Technology*, vol. 24, no. 1, pp. 177–186, 2013.
- [3] S. Xiong, M. Wei, M. Zhao, H. Xiong, W. Wang, and B. Zhou, "Hyperbolic tangent function weighted optimal intercept angle guidance law," *Aerospace Science and Technology*, vol. 78, pp. 604–619, 2018.
- [4] Z. Li and Z. Ding, "Robust Cooperative Guidance Law for Simultaneous Arrival," *IEEE Transactions on Control Systems Technology*, vol. 27, no. 3, pp. 1360–1367, May 2019.
- [5] S. He, H.-S. Shin, and A. Tsourdos, "Computational missile guidance: a deep reinforcement learning approach," *Journal of Aerospace Information Systems*, vol. 18, no. 8, pp. 571–582, 2021.
- [6] A. Ratnoo and D. Ghose, "Impact angle constrained interception of stationary targets," *Journal of Guidance, Control, and Dynamics*, vol. 31, no. 6, pp. 1817–1822, 2008.
- [7] I.-S. Jeon, J.-I. Lee, and M.-J. Tahk, "Impact-time-control guidance law for anti-ship missiles," *IEEE Transactions on control systems technology*, vol. 14, no. 2, pp. 260–266, 2006.
- [8] B. Gaudet, R. Furfaro, and R. Linares, "Reinforcement learning for angle-only intercept guidance of maneuvering targets," *Aerospace Science and Technology*, vol. 99, p. 105746, 2020.
- [9] H. M. La, R. Lim, and W. Sheng, "Multirobot cooperative learning for predator avoidance," *IEEE Transactions on Control Systems Technology*, vol. 23, no. 1, pp. 52–63, 2014.
- [10] H. Dong and X. Zhao, "Composite experience replay-based deep reinforcement learning with application in wind farm control," *IEEE Transactions on Control Systems Technology*, vol. 30, no. 3, pp. 1281–1295, 2021.
- [11] W. Kong, D. Zhou, Z. Yang, K. Zhang, and L. Zeng, "Maneuver strategy generation of ucav for within visual range air combat based on multi-agent reinforcement learning and target position prediction," *Applied Sciences*, vol. 10, no. 15, p. 5198, 2020.
- [12] B. M. Albaba and Y. Yildiz, "Driver modeling through deep reinforcement learning and behavioral game theory," *IEEE Transactions on Control Systems Technology*, vol. 30, no. 2, pp. 885–892, 2021.
- [13] T. Chen, K. Zhang, G. B. Giannakis, and T. Başar, "Communication-efficient policy gradient methods for distributed reinforcement learning," *IEEE Transactions on Control of Network Systems*, vol. 9, no. 2, pp. 917–929, 2021.
- [14] H. Dong, X. Zhao, and H. Yang, "Reinforcement learning-based approximate optimal control for attitude reorientation under state constraints," *IEEE Transactions on Control Systems Technology*, vol. 29, no. 4, pp. 1664–1673, 2020.
- [15] L. Chen, W. Wei hong, and L. Chao, "Deep reinforcement meta-learning guidance with impact angle constraint," *Journal of Astronautics*, 2021.
- [16] H. Zhang, H. Jiang, Y. Luo, and G. Xiao, "Data-driven optimal consensus control for discrete-time multi-agent systems with unknown dynamics using reinforcement learning method," *IEEE Transactions on Industrial Electronics*, vol. 64, no. 5, pp. 4091–4100, 2016.
- [17] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, "Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications," *IEEE transactions on cybernetics*, vol. 50, no. 9, pp. 3826–3839, 2020.
- [18] D. Brockhoff, A. Auger, N. Hansen, D. V. Arnold, and T. Hohm, "Mirrored sampling and sequential selection for evolution strategies," in *International Conference on Parallel Problem Solving from Nature*. Springer, 2010, pp. 11–21.
- [19] T. Salimans, J. Ho, X. Chen, S. Sidor, and I. Sutskever, "Evolution strategies as a scalable alternative to reinforcement learning," *arXiv preprint arXiv:1703.03864*, 2017.
- [20] D. Wierstra, T. Schaul, T. Glasmachers, Y. Sun, J. Peters, and J. Schmidhuber, "Natural evolution strategies," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 949–980, 2014.
- [21] B. Xu, Y. Zhang, D. Gong, Y. Guo, and M. Rong, "Environment sensitivity-based cooperative co-evolutionary algorithms for dynamic multi-objective optimization," *IEEE/ACM transactions on computational biology and bioinformatics*, vol. 15, no. 6, pp. 1877–1890, 2017.
- [22] H. Qu, K. Xing, and T. Alexander, "An improved genetic algorithm with co-evolutionary strategy for global path planning of multiple mobile robots," *Neurocomputing*, vol. 120, pp. 509–517, 2013.
- [23] Z. Wang, C. Chen, and D. Dong, "Instance weighted incremental evolution strategies for reinforcement learning in dynamic environments," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [24] P. Larrañaga and J. A. Lozano, *Estimation of distribution algorithms: A new tool for evolutionary computation*. Springer Science & Business Media, 2001, vol. 2.
- [25] H. Karshenas, R. Santana, C. Bielza, and P. Larrañaga, "Regularized continuous estimation of distribution algorithms," *Applied Soft Computing*, vol. 13, no. 5, pp. 2412–2432, 2013.
- [26] H. M. Omar and M. Abido, "Multiobjective evolutionary algorithm for designing fuzzy-based missile guidance laws," *Journal of Aerospace Engineering*, vol. 24, no. 1, pp. 89–94, 2011.
- [27] N. Maheswaranathan, L. Metz, G. Tucker, D. Choi, and J. Sohl-Dickstein, "Guided evolutionary strategies: Augmenting random search with surrogate gradients," in *International Conference on Machine Learning*. PMLR, 2019, pp. 4264–4273.
- [28] J. Del Ser, E. Osaba, D. Molina, X.-S. Yang, S. Salcedo-Sanz, D. Camacho, S. Das, P. N. Suganthan, C. A. C. Coello, and F. Herrera, "Bio-inspired computation: Where we stand and what's next," *Swarm and Evolutionary Computation*, vol. 48, pp. 220–250, 2019.
- [29] R. Gray, A. Franci, V. Srivastava, and N. E. Leonard, "Multiagent decision-making dynamics inspired by honeybees," *IEEE Transactions on Control of Network Systems*, vol. 5, no. 2, pp. 793–806, 2018.
- [30] B. Gaudet, R. Furfaro, and R. Linares, "Reinforcement Learning for Angle-Only Intercept Guidance of Maneuvering Targets," vol. 99, p. 105746. [Online]. Available: <http://arxiv.org/abs/1906.02113>
- [31] D. Han and S. Balakrishnan, "State-constrained agile missile control with adaptive-critic-based neural networks," *IEEE Transactions on Control Systems Technology*, vol. 10, no. 4, pp. 481–489, 2002.
- [32] H. Li, Y. Wu, and M. Chen, "Adaptive fault-tolerant tracking control for discrete-time multiagent systems via reinforcement learning algorithm," *IEEE Transactions on Cybernetics*, vol. 51, no. 3, pp. 1163–1174, 2020.
- [33] B. Thananjeyan, A. Balakrishna, S. Nair, M. Luo, K. Srinivasan, M. Hwang, J. E. Gonzalez, J. Ibarz, C. Finn, and K. Goldberg, "Recovery rl: Safe reinforcement learning with learned recovery zones," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4915–4922, 2021.
- [34] K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi, "Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning," in *International conference on machine learning*. PMLR, 2019, pp. 5887–5896.
- [35] Y.-J. Gong, W.-N. Chen, Z.-H. Zhan, J. Zhang, Y. Li, Q. Zhang, and J.-J. Li, "Distributed evolutionary algorithms and their models: A survey of the state-of-the-art," *Applied Soft Computing*, vol. 34, pp. 286–300, 2015.
- [36] A. Mendiburu, J. A. Lozano, and J. Miguel-Alonso, "Parallel implementation of edas based on probabilistic graphical models," *IEEE Transactions on Evolutionary Computation*, vol. 9, no. 4, pp. 406–423, 2005.
- [37] S. Trenn, "Multilayer perceptrons: Approximation order and necessary number of hidden units," *IEEE transactions on neural networks*, vol. 19, no. 5, pp. 836–844, 2008.
- [38] T. Lyu, Y. Guo, C. Li, G. Ma, and H. Zhang, "Multiple missiles cooperative guidance with simultaneous attack requirement under directed topologies," *Aerospace Science and Technology*, vol. 89, pp. 100–110, 2019.