

# 面向多机协同的 Att-MADDPG 围捕控制方法设计

刘 峰, 魏瑞轩, 丁 超, 姜龙亭, 李 天  
(空军工程大学航空工程学院, 西安, 710051)

**摘要** 多无人机对动态目标的围捕是无人机集群作战中的重要问题。针对面向动态目标的集群围捕问题, 通过分析基于 MADDPG 算法的围捕机制的不足, 借鉴 Google 机器翻译团队使用的注意力机制, 将注意力机制引入围捕过程, 设计基于注意力机制的协同围捕策略, 构建了相应的围捕算法。基于 AC 框架对 MADDPG 进行改进, 首先, 在 Critic 网络加入 Attention 模块, 依据不同注意力权重对所有围捕无人机进行信息处理; 然后, 在 Actor 网络加入 Attention 模块, 促使其他无人机进行协同围捕。仿真实验表明, Att-MADDPG 算法较 MADDPG 算法的训练稳定性提高 8.9%, 任务完成耗时减少 19.12%, 经学习后的围捕无人机通过协作配合使集群涌现出更具智能化围捕行为。

**关键词** 协同围捕; 强化学习; MADDPG; 智能性涌现

**DOI** 10.3969/j.issn.1009-3516.2021.03.002

**中图分类号** TP391.9    **文献标志码** A    **文章编号** 1009-3516(2021)03-0009-06

## Design of Att-MADDPG Hunting Control Method for Multi-UAV Cooperation

LIU Feng, WEI Ruixuan, DING Chao, JIANG Longting, LI Tian  
(Aeronautical Engineering College, Air Force Engineering University, Xi'an 710051, China)

**Abstract** The hunting of dynamic targets by multi-UAV is an important problem in UAV swarm operations. In this paper, aiming at the dynamic target oriented swarm hunting problem, by analyzing the shortcomings of the hunting mechanism based on MADDPG algorithm, and learning from the attention mechanism used by Google machine translation team, we introduce the attention mechanism into the hunting process, design the cooperative hunting strategy based on the attention mechanism, and construct the corresponding hunting algorithm. Improve MADDPG based on AC framework. First of all, the attention module is added to critical network to process the information of all UAVs according to different attention weights; then, the attention module is added to actor network to promote other UAVs to carry out cooperative hunting. The simulation results show that Att-MADDPG algorithm can improve the training stability by 8.9% and reduce the task completion time by 19.12% compared with MADDPG algorithm. After learning, the UAV can cooperate to make the swarm emerge more intelligent behavior.

**Key words** cooperative hunting; reinforcement learning; MADDPG; intelligence emergence

---

收稿日期: 2021-03-31

基金项目: 科技部“新一代人工智能”重点项目(2018AAA0102403)

作者简介: 刘 峰(1996—), 男, 河北邯郸人, 硕士生, 研究方向: 飞行器导航制导与智能控制。E-mail: 287321956@qq.com

通信作者: 魏瑞轩(1968—), 男, 陕西岐山人, 教授, 博士生导师, 研究方向: 智能飞行控制理论及应用。E-mail: ruixuanWeil23@163.com

**引用格式:** 刘峰, 魏瑞轩, 丁超, 等. 面向多机协同的 Att-MADDPG 围捕控制方法设计[J]. 空军工程大学学报(自然科学版), 2021, 22(3): 9-14. LIU Feng, WEI Ruixuan, DING Chao, et al. Design of Att-MADDPG Hunting Control Method for Multi-UAV Cooperation[J]. Journal of Air Force Engineering University (Natural Science Edition), 2021, 22(3): 9-14.

无人机集群是一种新型作战样式,具有作战成本低、冲突胜算大、生存能力强、作战效能高的特点。这些重要特征使得无人机集群在局部冲突中扮演着越来越重要的角色<sup>[1-2]</sup>。协同围捕问题属于无人机集群作战的典型应用场景之一,有重要的理论研究价值和广泛的应用前景。

在协同围捕方法设计上,已有不少学者做了相关研究<sup>[3-7]</sup>。张红强等<sup>[3]</sup>设计了一种基于简化虚拟受力模型,借助势域函数使机器人在未知动态环境下完成围捕。李瑞珍等<sup>[4]</sup>采用协商法为机器人分配动态围捕点,建立包含围捕路径损耗和包围效果的目标函数并优化航向角,从而实现协同围捕。Michael Rubenstein 等<sup>[5]</sup>以 1 000 个机器人为载体,分边缘检测、梯度上升、协同定位 3 部分算法设计,并通过局部交互进行合作,完成给定图片的不规则图形围捕演示,以人工集群的手段汇聚出自然蜂群的能力。

以上研究均是基于分布式控制,将协同围捕问题转换为集群任务分配、路径规划、群体一致性问题,从而达到围捕的效果,但在群体智能涌现方面仍有待提升。

近年来,也有部分学者探索通过强化学习方法来解决协同围捕问题<sup>[8-10]</sup>。吴子沉等<sup>[8]</sup>将围捕行为离散化后,设计能够应对复杂环境的围捕策略,但其存储机制仍有待优化。陈亮等<sup>[9]</sup>提出混合 DDPG 算法,有效协同异构 agent 之间的工作,同时,Q 函数重要信息丢失及过估计等问题有待解决。Ryan Lowe<sup>[10]</sup>于 2017 年提出 MADDPG 算法,采用“集中训练,分散执行”的框架解决了环境不稳定的问题,但是该算法随着 agent 数目的增加,Actor-Critic 网络难以训练和收敛。

针对以上分析,本文提出一种多无人机协同围捕算法 Att-MADDPG(即 Attention-MADDPG)。

## 1 问题描述

### 1.1 围捕环境描述

在一个无限大且无障碍的二维环境中,随机分布  $n(n \geq 3)$  架围捕无人机  $U_i$  和一架目标无人机  $T$ ,其速度分别为  $v_i, v_T$ ,且满足  $v_i > v_T$ ,航向角分别为  $\varphi_i, \varphi_T$ ,其中  $i \in I_n = \{1, 2, \dots, n\}$ ,如图 1 所示。

假设对于任意无人机,都能将自身的参数通过通信网络  $G$  实现与其他无人机实时信息交换。本文的目的就是基于这种信息共享,设计控制方法,使  $n$  架围捕无人机通过协作,在有限时间  $t$  内,在目标无人机周围形成围捕包围圈,迫使目标无人机  $T$  停止运动,从而完成围捕任务。理想的围捕包围圈通常是围捕无人机群均匀分布在以目标无人机  $T$  为圆心,围捕半径  $r$  的圆上<sup>[5]</sup>,以  $n=4$  为例,理想的围捕包围圈如图 2 所示。

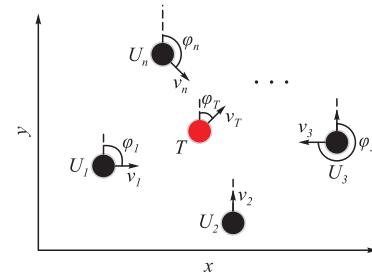


图 1 围捕环境示意图

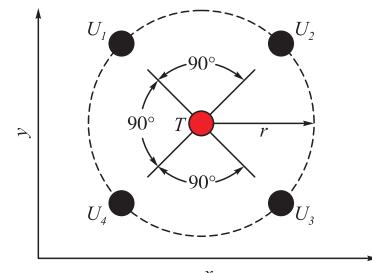


图 2 围捕包围圈示意图

### 1.2 围捕无人机模型

设无人机  $i$  当前时刻的位置为  $[x, y]^T$ ,构建非线性运动学方程如下:

$$\begin{cases} \dot{x}_i = v_i \cos(\varphi_i) \\ \dot{y}_i = v_i \sin(\varphi_i) \\ \dot{\varphi}_i = u_i \end{cases} \quad (1)$$

式中:  $u_i$  为无人机的控制输入,  $u_i \in [-\omega_0, \omega_0]$ ;  $\omega_0$  为无人机角速度上限。围捕控制策略就是根据围捕态势确定每架无人机的  $u_i$ ,使围捕无人机集群实现对目标无人机的有效围捕。

## 2 基于强化学习的 MADDPG 算法

### 2.1 强化学习理论

强化学习(reinforce learning, RL)是机器学习的一种,不同于监督学习或无监督学习,强化学习是通过与环境的交互进行不断试错-学习,以达成回报最大化或实现特定目标,如图 3 所示。

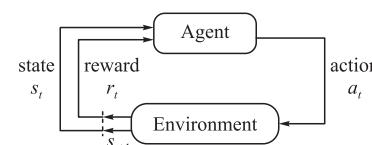


图 3 智能体与环境交互图

强化学习的常见模型是标准的马尔可夫决策过程(markov decision process, MDP)。由四元组  $(S, A, R, P_{s,a})$  表示,其中,  $S$  表示状态集,  $A$  表示动作

集,  $R$  表示奖励函数,  $P_{s,a}$  表示状态转移概率。基于当前状态  $s_t$ , 执行动作  $a_t$ , 以一定的状态转移概率达到下一时刻状态  $s_{t+1}$ , 获得即时奖励  $R_t$ , 但强化学习是寻找最大化累积回报的学习过程<sup>[11]</sup>。定义累积奖励期望值  $Q^\pi(s,a)$ :

$$Q^\pi(s,a) = E\left(\sum_{k=0}^{\infty} \gamma^k R_{t+k}\right) \quad (2)$$

式中:  $\gamma$  为折扣因子,  $0 < \gamma \leq 1$ , 表示注重长期奖励的程度。 $\pi$  为策略, 即状态到动作的映射。给出 Q-Learning 算法中  $Q$  值迭代计算表达式:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_t + \gamma \max Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$$

且  $s_t$  状态下最优策略为:

$$\pi^*(s) = \arg \max_{a \in A} Q(s, a)$$

Dietterich. T. G<sup>[12]</sup> 从值函数分解的角度, 完成了 Q-Learning 算法中  $Q$  值累加的收敛性证明。

## 2.2 MADDPG 算法

多智能体深度确定性策略梯度算法 (multi-agent deep deterministic policy gradient, MADDPG) 是对深度确定性策略梯度 (deep deterministic policy gradient, DDPG) 算法进行拓展, 使其能够适用于传统强化方法无法处理的多智能体合作问题的一种智能算法<sup>[10]</sup>。MADDPG 算法采用“集中训练, 分散执行”的框架进行学习, 见图 4。

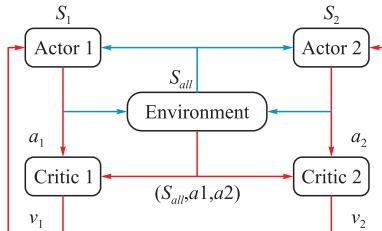


图 4 MADDPG 算法流程图

每个 agent 拥有一套独立的 Actor-Critic 网络, 其中 Actor 为行动网络, Critic 为评价网络。在训练过程中(图 4 中红色部分), 每个 agent 的 Critic 获取全信息状态, 同时包含所有 agent 的动作。当模型训练完毕, 每个 agent 的 Actor 利用局部信息完成与环境的交互(图 4 中蓝色部分)。令  $O_i$  表示 agent 对环境的观测,  $I_A$  和  $a_i$  分别表示 Actor 的输入和输出,  $I_C$  和  $Q_i$  分别表示 Critic 的输入和输出, 那么对于第  $i$  个 agent, 网络输入输出为:

$$\begin{aligned} I_A &= O_i, a_i = \pi_i(s_i) \\ I_C &= [O_i, a_1, a_2, \dots, a_n], Q_i = Q_i(I_C) \end{aligned} \quad (3)$$

当 agent 数量增多时, 由式(3)中  $I_C$  可知, Critic 的输入维度也呈线性增长, 这将导致网络难以训练和收敛。文献[10]和[13]同样指出, 尽管集中训练, 分散执行的结构具有诸多优势, 但是随着 agent 数量的增加, 集中训练中 Critic 网络规模会快速增长,

因而无法处理大规模多智能体的学习问题。同样的, 由 Facebook AI 实验室和 Google AI 联合赞助的二维网格环境炸弹人平台, 在测试时最多也只能容纳 4 个 agent。

## 3 基于 Att-MADDPG 的围捕控制策略设计

### 3.1 面向无人机围捕的 Attention 机制

近年来, 注意力机制 (attention mechanism) 被广泛用于基于深度学习的自然语言处理、图像分类、机器翻译可视化对齐、语音识别等各种任务中, 并取得了不错的效果<sup>[14]</sup>。2017 年 6 月, Google 机器翻译团队借助自注意力机制在 WMT2014 语料中的英德和英法翻译任务上取得了优异成绩, 翻译错误率降低了 60%, 并且训练速度远优于其他主流模型<sup>[15]</sup>。

围捕过程中也存在类似的注意力问题。每架无人机更多的关注与自己近邻的无人机, 对距离较远的无人机的态势关注的较少, 甚至不关注。这就是围捕过程中的注意力现象。我们将这种现象引入到围捕策略的设计, 形成面向围捕的注意力机制。以无人机协同围捕的场景程阐释注意力机制如下:

定义围捕无人机集群的联合动作为  $Source$ , 待处理信息为  $Target$ :

$$Source = \langle p_1, p_2, \dots, p_m \rangle$$

$$Target = \langle q_1, q_2, \dots, q_n \rangle$$

其中,  $p_i (i=1, 2, \dots, m)$  表示第  $i$  架无人机动作,  $q_i (i=1, 2, \dots, n)$  表示待处理信息。Attention 机制<sup>[16-17]</sup> 最常用的是编码器-解码器 (Encoder-Decoder) 框架, 如图 5 所示。Encoder 对输入的  $Source$  进行编码, 通过神经网络的非线性变换转化为注意力分布  $C$ ,  $C = \{c_1, c_2, \dots, c_{L_p}\}$ , 其中  $L_p$  为  $Source$  的长度, Decoder 根据注意力分布  $C$  和  $n-1$  时刻无人机的位置生成  $n$  时刻的信息  $q_n$ , 即围捕无人机集群待处理信息。给出注意力分布  $c_i (i=1, 2, \dots, L_p)$  的表达式:

$$\begin{aligned} c_i &= \sum_{j=1}^{L_p} w_{ij} p_j \\ q_n &= \text{Attention}(c_i, q_{n-1}) \end{aligned} \quad (4)$$

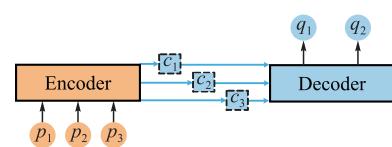


图 5 基于 Attention 机制的 Encoder-Decoder 框架  
式中:  $w_{ij}$  为  $Source$  中第  $j$  架无人机的注意力权重系

数;  $p_j$  为 Source 中第  $j$  架无人机的动作信息; Attention 为非线性变换函数。

给出基于 Attention 机制下注意力分布  $c_i$  的具体计算过程:

1) 计算 2 架围捕无人机之间的相关性系数:  
 $Similarity_i = \ln(Distance_{ij}/D)$ , 其中  $Distance_{ij}$  为两架围捕无人机之间距离,  $D$  为有效利用区域半径。

2) 引入 Softmax 函数对第 1 阶段的相关性系数进行归一化处理, 得到注意力权重系数  $\omega_i$ 。一方面将原始分值映射成所有元素权重之和为 1 的概率分布, 另一方面通过 Softmax 的内在机制突出重要元素的注意力权重系数。

$$\omega_i = \text{Softmax}(Similarity_i) = \frac{e^{Similarity_i}}{\sum_{j=1}^{L_p} e^{Similarity_j}} \quad (5)$$

3) 根据注意力权重系数对围捕无人机信息进行加权求和, 计算注意力分布  $c_i$  值。

### 3.2 基于 Att-MADDPG 的围捕控制策略

本文 2.2 节中分析 MADDPG 算法由于  $I_C$  的计算中使用了所有 agent 的信息, 使得其训练收敛受到影响。因此, 我们引入 Attention 机制对信息进行注意力筛选, 从而提高信息的有效利用率, 算法框图如图 6 所示。

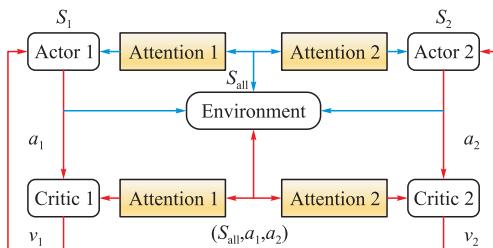


图 6 基于 Attention 机制改进的 MADDPG 算法框图

与 MADDPG 算法不同之处在于中心化的大脑(即 Critic)协调所有围捕无人机的动作之前, 经各自 Attention 模块进行非线性处理, 对有效利用区域内的围捕无人机信息进行策略评估(图 6 中红色部分所示)。当模型训练完毕, 依据 Actor 利用局部信息完成与环境的交互(图 6 中蓝色部分所示)。则围捕无人机  $i$  的有效利用区域值函数为:

$$Q_i^\pi(s, a) = \text{Attention}[O_i, a_1, a_2, \dots, a_n]$$

因此, 每架围捕无人机的 Critic 网络拟合的是有效利用区域的全局值函数, 而非围捕无人机自身的值函数。这样, 只需要围捕无人机的策略朝着有效利用区域的全局值函数的方向更新即可。使用 MADDPG 算法中双网络进行更新:

$$\begin{aligned} L(\theta_i) &= E[(y_i - Q_i^\pi(s, a))^2] \\ y_i &= R + \gamma Q_i^\pi(s', a') \mid_{a_i = \pi'_{\theta_i}(s_i)} \end{aligned} \quad (6)$$

式中:  $y_i$  为目标网络的值函数, 由即时奖励和下一步确定策略值函数构成;  $L(\theta_i)$  为目标 Critic 网络损失函数,  $\theta_i$  为网络中参数集合。目标 Actor 网络和目标 Critic 网络采用周期性平稳滑动方法从 Actor-Critic 网络中复制参数进行更新。目标 Critic 网络损失梯度通过链式法则进行求导, 其梯度为:

$$\nabla L(\theta_i) = E[\nabla_{a_i} Q_i^\pi(s, a) \nabla_{\theta_i} \pi'_{\theta_i}(s_i) \mid a_i = \pi'_{\theta_i}(s_i)] \quad (7)$$

式中: 状态  $s$  为有效利用区域全局观测;  $s_i$  为围捕无人机自身观测。Att-MADDPG 算法流程见图 7。

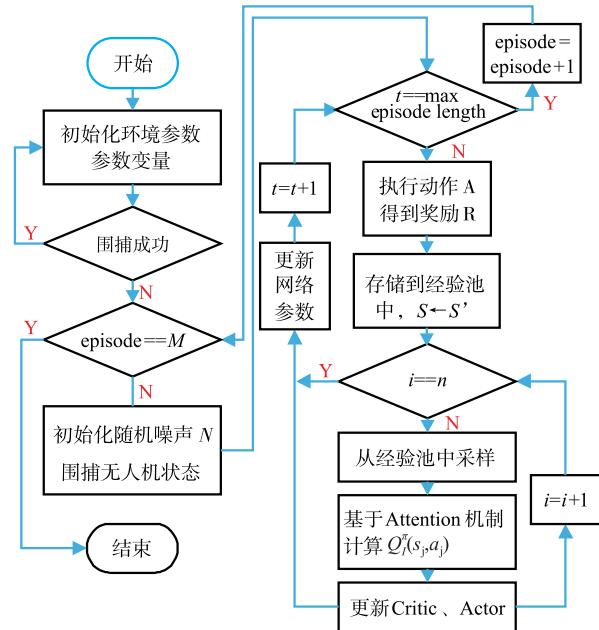


图 7 Att-MADDPG 算法流程图

Att-MADDPG 算法伪代码如下所示:

```

初始化环境参数、参数变量;
判断是否围捕成功;
for episode=1 to M do
    初始化随机噪声 N;
    初始化围捕无人机状态 S ;
    for t=1 to max-episode-length do
        每架围捕无人机采用随机策略执行一次动作 A ,与环境交互后得到即时奖励 R ,并达到新的状态 S' ;
        存储 (S,A,R,S') 到经验池中, S ← S' ;
        for 围捕无人机 i=1 to n do
            从经验池中采样 (S^i, A^i, R^i, S'^i) ;
            基于 Attention 机制计算 Q_i^pi(S^i, A^i) ;
            设置目标网络值函数 y^i ;
            基于最小损失函数 L(theta_i) 更新 Critic;
            基于策略梯度 grad L(theta_i) 更新 Actor;
        end for
        更新每架围捕无人机网络参数:
        theta'_i ← alpha * theta_i + (1 - alpha) * theta'_i;
    end for
end for

```

## 4 仿真实验

为验证所设计 Att-MADDPG 算法的有效性及智能性,取围捕无人机数量  $n=4$  进行动态协同围捕仿真实验,并对比 MADDPG 算法进行训练,测试相关性能指标。

### 4.1 仿真环境配置

设置围捕无人机奖励函数如下:

$$R^i = \underbrace{100 \times (Dis_{t-1}^i - Dis_t^i)}_{\textcircled{1}} - \underbrace{t/1000}_{\textcircled{2}} + \underbrace{10 \times flag}_{\textcircled{3}}$$

$$flag = \begin{cases} 1, & \text{形成围捕包围圈} \\ 0, & \text{未形成围捕包围圈} \end{cases}$$

式中: $Dis_t^i$  为  $t$  时刻围捕无人机  $i$  与目标无人机 T 的距离;flag 为标志位。奖励函数由 3 部分构成,第①部分引导围捕无人机靠近目标无人机,整体乘以 100,突出动作的奖励,第②部分为围捕过程中所消耗的时间,第③部分促使围捕无人机集群形成围捕包围圈。

仿真环境参数设置如表 1 所示。

表 1 仿真场景设置

参数物理含义	参数数值
目标无人机速度 $v_T$ / (km/h)	250
围捕无人机 1 速度 $v_1$ / (km/h)	350
围捕无人机 2 速度 $v_2$ / (km/h)	300
围捕无人机 3 速度 $v_3$ / (km/h)	300
围捕无人机 4 速度 $v_4$ / (km/h)	300
有限时间 $t$ / s	300
围捕半径 $r$ / km	1.5
角速度上限 $\omega_0$ / (rad/s)	$\frac{\pi}{6}$
有效利用区域半径 $D$ / km	6

配置相同的四机协同围捕环境,设置相同的奖励函数,MADDPG 算法及 Att-MADDPG 算法的训练过程如图 8~9 所示。如图 8 所示,采用 MADDPG 算法训练的围捕无人机集群大致完成协同围捕任务,但训练过程中仍存在一些片段无法有效围捕。

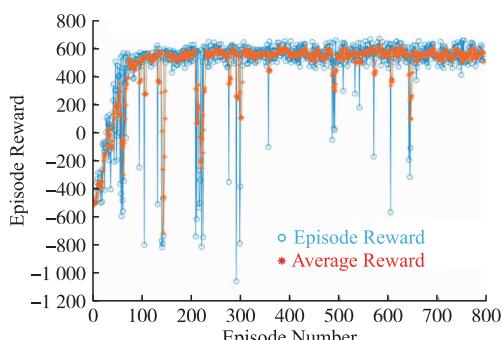


图 8 MADDPG 算法训练过程图

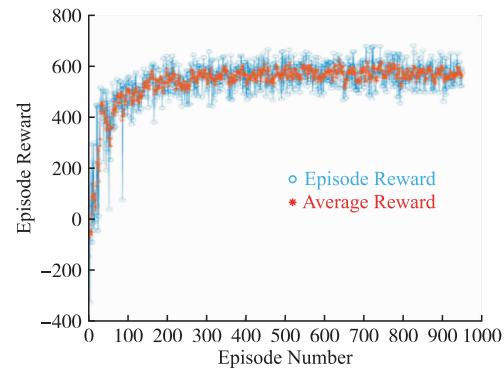


图 9 Att-MADDPG 算法训练过程图

由图 9 可知,经过训练,Att-MADDPG 算法的平均奖励在 150 片段后收敛,并大致均匀分布 570 左右,围捕无人机集群能够有效完成协同围捕任务,并获得较高奖励。

引入 Att-MADDPG 算法及 MADDPG 算法在稳定收敛后 1 000 个片段的均方差,对比如验证算法的稳定性,见表 2。Att-MADDPG 算法较 MADDPG 算法在协同围捕上更加稳定,根据  $3\sigma$  原则,记录稳定收敛后 1 000 个片段不在此范围内的平均奖励数目,见表 3,Att-MADDPG 算法稳定性较 MADDPG 算法提高 8.9%。

表 2 稳定收敛后 1 000 个片段的均方差

MADDPG 算法	Att-MADDPG 算法
107.679 2	25.897 8

表 3 1 000 个片段内不在  $3\sigma$  原则范围内的数目

MADDPG 算法	Att-MADDPG 算法
102	13

### 4.2 动态协同围捕仿真验证

设定目标无人机的运动策略为固定直线轨迹,对经过 Att-MADDPG 算法训练后的 4 架智能围捕无人机进行验证,验证集参数见表 4。

表 4 动态协同围捕验证集参数

围捕无人机	初始位置及航向
$U_1$	$[0, 0, 0]^T$
$U_2$	$[1, 8, \pi/2]^T$
$U_3$	$[13, 1, \pi/6]^T$
$U_4$	$[3, 4, \pi/2]^T$

由图 10 可知,围捕无人机通过相互协作完成围捕。在围捕过程中,无人机实时判断围捕态势,引入注意力机制观察有效利用区域半径内的其他无人机状态信息,达到形成围捕包围圈的目的,各无人机经过协作使系统涌现出更加智能化的协同围捕行为。

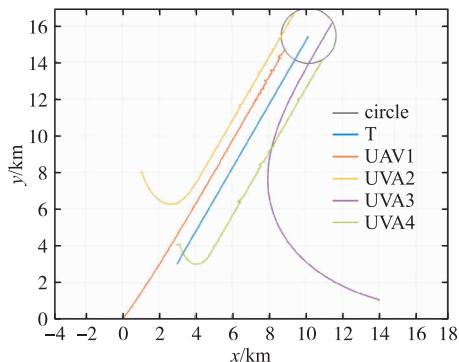


图 10 动态协同围捕轨迹图

环境配置相同,采用 MADDPG 算法进行验证,Att-MADDPG 算法完成协同围捕总用时 264 s,MADDPG 算法本文算法总用时 326.4 s,减少 19.12%。

## 5 结论

无人机集群作战在局部冲突中发挥着越来越重要的作用,协同围捕是无人机集群作战的典型应用场景之一,也是集群作战中的重要问题。本文针对 MADDPG 算法随着 agent 数量的增加,训练难以收敛的不足,基于注意力机制提出 Att-MADDPG 围捕控制方法,较 MADDPG 算法的训练稳定性提高 8.9%,任务完成耗时减少 19.12%,且经学习后的围捕无人机通过协作配合使集群涌现出更具智能化围捕行为。

为使本文所提算法能够适用于更加复杂的环境,仍需研究基于群智汇聚的协同围捕机理,并优化在三维环境下的协同围捕策略,使围捕行为更具智能化。

## 参考文献

- [1] 魏瑞轩,李学仁. 无人机系统及作战使用[M]. 北京: 国防工业出版社, 2009.
- [2] 郝雅楠,孔超,关晓红. 国外无人机作战运用与发展趋势分析:关于纳卡冲突事件的思考[J]. 国防科技工业, 2021(2):50-53.
- [3] 张红强,吴亮红,周游,等. 复杂环境下群机器人自组织协同多目标围捕[J]. 控制理论与应用, 2020, 37(5):1054-1062.
- [4] 李瑞珍,杨惠珍,萧丛杉. 基于动态围捕点的多机器人协同策略[J]. 控制工程, 2019, 26(3):510-514.
- [5] MICHAEL R, ALEJANDRO , RADHIKA N. Programmable Self-Assembly in a Thousand-Robot Swarm [J]. Science, 345(6198), 795-799.
- [6] 徐雪松,曾智,邵红燕,等. 基于个体-协同触发强化学习的多机器人行为决策方法[J]. 仪器仪表学报, 2020, 41(5):66-75.
- [7] 徐雪松,杨胜杰,陈荣元. 复杂环境移动群机器人最优路径规划方法[J]. 电子测量与仪器学报, 2016, 30(2):274-282.
- [8] 吴子沉,胡斌. 基于态势认知的无人机集群围捕方法[J]. 北京航空航天大学学报, 2021, 47(2):424-430.
- [9] 陈亮,梁宸,张景异,等. Actor-Critic 框架下一种基于改进 DDPG 的多智能体强化学习算法[J]. 控制与决策, 2021, 36(1): 75-82.
- [10] LOWE R, WU Y I, TAMAR A, et al. Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments[C]// Advances in Neural Information Processing Systems. 2017:6379-6390.
- [11] BIANCHI R A C, RIBEIRO C H C, COSTA A H R. Accelerating autonomous learning by using heuristic selection of actions [J]. Journal of Heuristics, 2008, 14(2):135-168.
- [12] DIETTERICH T G. Hierarchical Reinforcement Learning With the MAXQ Value Function Decomposition [J]. Journal of Artificial Intelligence Research, 2000(13):227-303.
- [13] 孙彧,曹雷,陈希亮,等. 多智能体深度强化学习研究综述[J]. 计算机工程与应用, 2020, 56(5):13-24.
- [14] YIN W, EBERT S, SCHÜTZE H. Attention-Based Convolutional Neural Network for Machine Comprehension[J]. Computer Science, 2016(7):15-21.
- [15] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is All You Need[C]//31st Conference on Neural Information Processing Systems. 2017: Long Beach, [s. n.]:1-15.
- [16] SHARIQ I, FEI S. Actor-Attention-Critic for Multi-Agent Reinforcement Learning[C]//Proceedings of the 36th International Conference on Machine Learning, PMLR 97. 2019:2961-2970,
- [17] CHEN C, LIU Y, KREISS S, ALAHI A. Crowd-Robot Interaction: Crowd-Aware Robot Navigation With Attention-Based Deep Reinforcement Learning [C]// 2019 International Conference on Robotics and Automation (ICRA), 2019: 6015-6022.

(编辑:徐敏)