

基于深度强化学习的短波自主检测与定位方法

冯祺玥, 唐涛, 张昀普, 王鼎, 吴志东
(信息工程大学信息工程学院, 河南 郑州 450001)

摘要: 在短波信号实时检测与定位环节中, 往往存在数据关联与误差传递的问题, 易导致检测与定位信号不匹配, 且难以实时输出结果。针对上述问题, 提出一种基于多智能体近端策略优化(Multi-agent proximal policy optimization, MAPPO)的短波自主检测与定位方法 MAPPO-DL (Detection and localization based on MAPPO)。通过设计基于短波信号时频图的强化学习环境, 有效表征短波信号的频域特征。同时, 设计具有混合动作空间的滤波窗智能体以自适应选择信号。另外, 通过定位误差椭圆概率、信号匹配度与多智能体协作策略设计奖励函数, 利用改进的 MAPPO 算法探索最优策略, 实现短波自主检测与定位。经过仿真测试, 相比于基线算法, 信号检测与定位平均时间缩短了 0.12s, 平均检测与定位准确率提高了 22%, 为复杂电磁环境下自主协同目标检测与定位提供新思路。

关键词: 短波信号检测与定位; 多智能体近端策略优化; 强化学习; 混合动作空间; 定位误差椭圆概率

中图分类号: V247; TN97

文献标志码: A

DOI: 10.11959/j.issn.1000

A Shortwave Autonomous Signal Detection and Localization Method Based on Deep Reinforcement Learning

Feng Qiyue, Tang Tao, Zhang Yunpu, Wang Ding, Wu Zhidong

School of Information System Engineering, Information Engineering University, Zhengzhou 450001, China

Abstract: In the process of real-time shortwave signal detection and localization, problems such as data association and error propagation were frequently observed, by which a mismatch between detection and localization results was easily induced, and the timely output of outcomes was impeded. To address the above issues, an integrated shortwave signal detection and localization algorithm based on Multi-agent Proximal Policy Optimization (MAPPO-DL) was proposed. By constructing a reinforcement learning environment using shortwave signal time-frequency diagrams, spectral characteristics of shortwave signals were effectively captured. A filtering-window agent with hybrid action space was designed for adaptive signal selection. Furthermore, localization error ellipse probability, signal matching degree, and multi-agent collaboration strategies were incorporated into the reward function design. Optimal policies were explored by the optimized MAPPO to achieve autonomous shortwave signal detection and localization. Simulation results demonstrate that compared with baseline algorithms, our approach reduces average recognition-localization time by 0.12s while improving accuracy by 22%, thereby providing a novel solution for autonomous cooperative target detection and localization in complex electromagnetic environments.

Keywords: Shortwave signal detection and localization, Multi-Agent Proximal Policy Optimization, Reinforcement learning, Hybrid action space, Localization error ellipses

收稿日期: XXXX-XX-XX; 修回日期: XXXX-XX-XX

通信作者: 唐涛, 13703820621@163.com

基金项目: 国家自然科学基金资助项目(No.62171469)

Foundation Items: The National Natural Science Foundation of China (No.62171469)

0 引言

短波通信凭借其独特的超视距传输能力,在应急通信、军事侦察及远洋导航等领域具有不可替代的优势,但受限于多径效应、动态衰减等复杂信道特性,导致短波信号检测与定位较为困难^[1-2]。短波常用定位体制为多站测角定位,需要利用波达方向(Direction of Arrival, DOA)估计方法得到入射信号方位角和仰角,然后利用参数估计方法确定目标位置,该方法依赖精确的角度测量结果^[3-4]。因此,经典信号检测与定位方法(Detection and Localization, DL)往往需要先检测发现信号,再进行检测信号类型和定位等工作,该过程需要大量人工参与判别测定信号等操作。在实时变化的信号环境下,传统方法往往存在数据关联与误差传递的问题,导致定位实时性较差,难以满足短波信号实时检测与定位的需求。

随着人工智能技术的快速发展,强化学习(Reinforcement Learning, RL)在无线通信领域的应用为信号自主检测与定位提供了新思路^[5]。单智能体深度强化学习通过与环境交互自主优化策略,已在非动态场景的频谱检测与功率控制中取得进展。在信号自主检测与定位领域,强化学习通过模拟智能体与环境交互,自适应地学习图像信号特征并优化分类策略^[6]。Caicedo等^[7-9]提出了基于深度Q网络(Deep Q Network, DQN)的图形信号检测分类任务。Li等^[10]提出了利用DQN算法实现无线信号定位的方法,设计了由无线接收信号强度中提取信号位置的奖励机制,经过模型训练可实现无线信号定位。Paul等^[11]提出一种基于强化学习的DOA位置估计方法,再利用最小二乘法估计信号源位置,以获得高精度的定位结果。

与传统方法相比,强化学习在动态干扰的低信噪比环境下表现出更强的鲁棒性和自适应能力^[12]。基于RL的定位方法通过建立定位场景的马尔可夫观测过程,动态优化测向准确性,从而提升定位精度与效率,表现出良好的定位性能。然而,短波测向定位系统是一个多智能体协作系统,各站点根据局部观测信息实时协同决策,传统单智能体RL方法忽视了智能体间的协作,导致定位效率受限。

多智能体强化学习(Multi-agent reinforcement learning, MARL)作为强化学习的研究热点,已经应用于无人机、物联网等方面^[13]。MARL通过模拟

多个智能体之间的协作与竞争,能够有效解决复杂环境下的决策问题^[14]。其中,多智能体近端策略优化算法(Multi-agent proximal policy optimization, MAPPO)作为一种高效的MARL的算法,通过策略优化和分布式学习,为多智能体系统的协同决策提供了新的思路^[15]。MAPPO更适用于实时变化的环境下的多目标信号分析任务,可以优化信号接收任务,具有较强的自主性和适应性。Wang等^[16]提出了一种基于多智能体深度强化学习和策略优化的方法用于多源信号定位问题,并通过智能体之间的协作实现高精度定位。另外,基于值分解的强化学习补偿滤波的多智能体协同定位算法被用来解决非线性环境下的定位问题^[17]。Alagha等^[18]提出使用MAPPO算法解决目标定位问题,但该算法仅适用于二维平面的定位问题,尚不能解决三维空间中信号的定位问题。

综上,现有基于MARL算法的信号检测与定位方法仍存在以下不足。首先,缺乏泛化能力,部分算法仅适用于特定目标的定位,无法在时变的信号环境中定位目标;其次,短波信号先检测再定位的两步方法中处理环节较多,在信号实时检测与定位环节中可能存在检测与定位信号不匹配问题。另外,部分算法仅研究基于MARL算法的信号检测或定位的单一方案,尚未将二者整体考虑,因此无法同时实现信号的检测与定位。

针对上述问题,本文基于MAPPO算法,提出了一种短波信号自主检测与定位的新方法MAPPO-DL (Detection and localization based on MAPPO)。通过充分利用信号特征的时空分布特性,对阵列数据进行滤波处理得到时频图。然后,在时频图基础上构建多智能体强化学习马尔可夫决策过程,通过智能体滤波信号同时完成检测与定位工作,以避免检测与定位信号不匹配问题。最后,引入定位误差椭圆概率、信号匹配度和多智能体协作策略设计奖励函数,实现信号高效检测与实时定位一体化。

1 系统模型

短波信号检测与定位模型如图1所示,短波信号接收设备接收信号后,将入射信号时域数据做快速傅里叶变换(Fast Fourier Transform, FFT),得到入射信号的以时间为横坐标,以频率为纵坐标的时频图。在时频图上滤波检测信号,以获得信号数

据,再根据定位算法获得目标位置。另一方面,以时频图作为输入,利用信号检测算法得到信号类型。

将地球建模为半径为等效地球半径的球体,并在该球体上建立了一个忽略高度的地理坐标系。 NR 表示北极,目标源和站点位于球体表面,其坐标用纬度 ρ 和经度 ω 表示。假设地球表面存在 D 个待定位目标,该目标源辐射短波信号。 $\mathbf{u} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_D]^T$ 表示 D 个未知的目标源位置,第 d 个目标源的经纬度为 $\mathbf{u}_d = [\omega_d, \rho_d]$, $d \in \{1, 2, \dots, D\}$,第 n 个测向站的经纬度为 $\mathbf{s}_n = [\omega_n, \rho_n]$, $n \in \{1, 2, \dots, N\}$ 。 $\theta_n \in (-\pi, \pi]$ 表示信号源 \mathbf{u}_d 到站点 \mathbf{s}_n 的方位角,方位角在正北方向为零,顺时针方向为正。

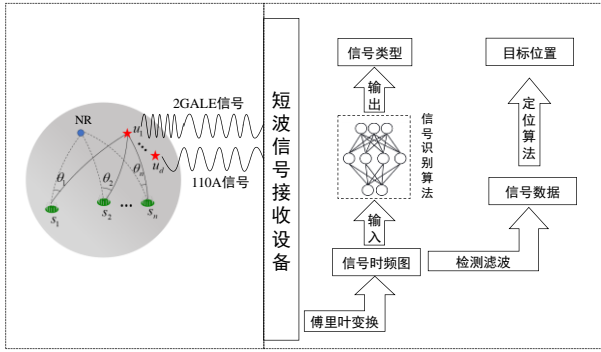


图1 短波信号检测与定位模型

每个测向站中包含 M 个阵元的阵列天线同步接收信号,设同时有 D 个具有相同中心频率 ω_0 ,波长为 λ 的空间窄带平面波($M > D$)分别以入射角 $\theta_1, \theta_2, \dots, \theta_D$ 入射该天线阵列,其中 θ_i 是第 i 个入射信号的方位角。窄带平面波信号的阵列输出数学模型矩阵形式为:

$$\mathbf{X}(t) = \mathbf{A}(\theta)\mathbf{S}(t) + \mathbf{N}(t) \quad (1)$$

其中, $\mathbf{X}(t) = [x_1(t), x_2(t), \dots, x_M(t)]$ 为阵列输出矢量, $\mathbf{N}(t) = [n_1(t), n_2(t), \dots, n_M(t)]$ 为阵列加性噪声矢量, $\mathbf{S}(t) = [s_1(t), s_2(t), \dots, s_D(t)]$ 为接收信源矢量, $\mathbf{A}(\theta) = [\mathbf{a}(\theta_1) \mathbf{a}(\theta_2) \dots \mathbf{a}(\theta_D)]$ 为阵列流型矩阵, $\mathbf{a}(\theta_i) = [e^{-j\omega_0\tau_{1i}}, e^{-j\omega_0\tau_{2i}}, \dots, e^{-j\omega_0\tau_{Mi}}]^T$ 为阵列方向矢量, ω_i 为信号 i 的中心频率。

对于 D 个目标源信号的到达角,MUSIC空间谱估计 $\hat{P}_{MUSIC}(\theta)$ 可以用噪声子空间 \mathbf{U}_M 的特征矢量 $\Phi_i, i = D + 1, \dots, M$ 来表示:

$$\hat{P}_{MUSIC}(\theta) = \mathbf{a}^H(\theta) \left(\sum_{i=D+1}^M \hat{\Phi}_i \hat{\Phi}_i^H \right) \mathbf{a}(\theta), i = D + 1, \dots, M \quad (2)$$

根据式(2)计算得到的空间谱采用谱峰搜索获得信号来向 $\theta_1, \theta_2, \dots, \theta_D$ 。根据方位角估计结果进行交叉定位,对所得交点取平均得到目标估计位置 $\hat{\mathbf{u}}_d$ 。

在信号检测任务中,以5类短波信号为例,分别是:2GALE、LINK4A、110A、CLOVER2000和LINK11。如图2所示,经过下变频采样设备接收的短波信号在时频图上的特点各有不同,2GALE、LINK4A采用FSK调制方式,其时频图呈现多条间隔出现的线段;110A采用PSK调制方式,时频图中有效区域呈矩形;CLOVER2000的时频图包含8条断断续续的斑点;LINK11的特殊帧结构导致其时频图由哑铃状的矩形组成。因此,通过时频图即可区分短波特定信号。

在动态变化的信号环境下,基于深度学习的智能检测算法一方面依赖大规模数据集进行训练,另一方面,信号间的相互干扰可能导致检测信号与定位信号不一致,同时也未能实现短波信号的测向与定位一体化。针对短波多信号动态性强、干扰大的问题,本文采用基于多智能体的深度强化学习对信号检测与定位问题进行MDP建模,使智能体能够在信号环境中自主学习和探索,从而实现短波信号的实时测向与定位。

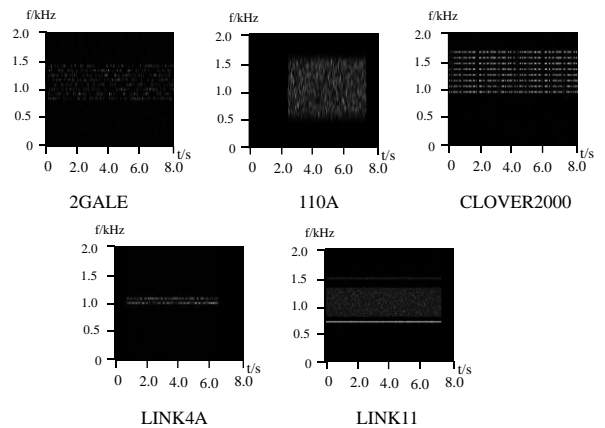


图2 五类短波信号时频图

2 MAPPO-DL

2.1 建立马尔可夫决策过程

本文提出的MAPPO-DL方案利用深度强化学

习机制，将短波信号检测与定位过程建模为可观测马尔可夫决策过程。智能体通过与环境的不断交互试错，以寻找累积奖励最大的策略，从而实现边执行边学习的自主进化能力。为实现短波信号自主检测和定位，需要先对入射信号做滤波处理得到空时频图，如图 3 所示。

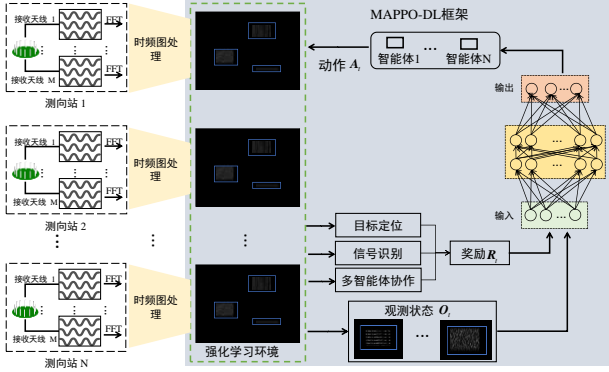


图 3 MAPPO-DL 框架

如图 3 所示，将入射信号时域数据做 FFT，得到入射信号的空间谱时频图。其次，对每个时间片数据按频域模型做 MUSIC 测向处理，记录每个时频点对应的测向结果，形成关于每个时间片数据的空间谱测向结果^[19]。在基于时频图的多智能体信号检测与定位环境中，多个测向站点同步接收信号。以下分别对该强化学习马尔可夫过程的动作、状态与奖励设计进行阐述。

2.1.1 动作

动作是指智能体根据当前状态为达到目标做出的决策。智能体是指做动作的主体，在 MAPPO-DL 环境中自适应滤波窗为智能体，智能体的长为 w ，宽为 h 。为使智能体能够快速准确选择信号并检测与定位，考虑设计离散动作 x_t^i 和连续动作 y_t^i 相

$$T_n = \begin{bmatrix} -\sin(\omega_n) & \cos(\omega_n) & 0 \\ -\cos(\omega_n)\sin(\rho_n) & -\sin(\omega_n)\sin(\rho_n) & \cos(\rho_n) \\ \cos(\omega_n)\cos(\rho_n) & \sin(\omega_n)\cos(\rho_n) & \sin(\rho_n) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_n^T \\ \mathbf{b}_n^T \\ \mathbf{c}_n^T \end{bmatrix} \quad (1 \leq n \leq N) \quad (6)$$

假设第 d 个短波信号源的经纬度为 (ω_d, ρ_d) ，于是该短波辐射源的地心地固坐标为：

$$\mathbf{u}_d = \begin{bmatrix} \gamma_d \cos(\rho_d) \cos(\omega_d) \\ \gamma_d \cos(\rho_d) \sin(\omega_d) \\ \gamma_d \sin(\rho_d) \end{bmatrix} \quad (7)$$

其中， $\gamma_d = R_e / \sqrt{1 - e^2 (\sin(\rho_d))^2}$ 。

结合的混合动作，则多智能体环境中的动作空间为：

$$\mathbf{A}_t = [a_t^1, a_t^2, \dots, a_t^n] \quad (3)$$

其中， $a_t^i = [x_t^i, y_t^i]^T$ ，离散动作 x_t^i 为滤波窗移动方向，主要包括上、下、左、右、左上、左下、右上和右下 8 个方向。需要说明的是左上、左下、右上和右下的方向均为 45° 方向；连续动作 $y_t^i = [w_t^i, h_t^i]$ 为自适应滤波窗的尺度变换。

2.1.2 状态

在信号智能检测与定位环境中， t 时刻的全局状态 s_t 为接收信号的时频图。智能体 i 的局部观测状态为其周围 $(2w_{\max}, 2h_{\max})$ 大小的图像 o_t^i ，使滤波窗智能体更好地检测执行动作后的环境变化。因此总体观测状态空间为：

$$\mathbf{O}_t = [o_t^1, o_t^2, \dots, o_t^n] \quad (4)$$

2.1.3 奖励

(1) 定位奖励：在短波空间谱测向过程中，测向的准确度与智能体所选信号的信噪比(Signal-to-noise, SNR) 呈正相关，即信噪比越高，测向精度越高；反之则越低。在任意多个观测站条件下同步测向交会定位时，首先将测向站经纬度转换为地心地固坐标。考虑利用 N 个测向站对地球表面的短波辐射源进行定位，第 n 个测向站的经纬度为 (ω_n, ρ_n) ，则该测向站的地心地固坐标为：

$$\mathbf{s}_n = \begin{bmatrix} \gamma_n \cos(\rho_n) \cos(\omega_n) \\ \gamma_n \cos(\rho_n) \sin(\omega_n) \\ \gamma_n \sin(\rho_n) \end{bmatrix} \quad (1 \leq n \leq N) \quad (5)$$

其中， $\gamma_n = R_e / \sqrt{1 - (\sin(\rho_n))^2}$ ， $R_e = 6378.160 \text{ km}$ 为地球等效半径。该测向站对应的坐标转换矩阵为：

根据各测向站对目标测得的测向角度 θ_i ，可以求得 N 条沿地球表面的曲线，交叉定位形成 C_N^2 个目标估计点 $\hat{\mathbf{u}}_d = [\hat{u}_d^1, \hat{u}_d^2, \dots, \hat{u}_d^{C_N^2}]$ 。假设目标估计误差服从零均值高斯分布，协方差矩阵为 $\mathbf{R}_{uu} = E[(\hat{\mathbf{u}}_d - \mathbf{u}_d)(\hat{\mathbf{u}}_d - \mathbf{u}_d)^T]$ ，则目标估计的概率密度函数为：

$$p_{\hat{u}_d}(\xi) = (2\pi)^{-n/2} (\det [R_{uu}])^{-1/2} \cdot \exp\left\{-\frac{(\xi - \hat{u}_d)^T R_{uu}^{-1} (\xi - \hat{u}_d)}{2}\right\} \quad (8)$$

其中, n 表示 \mathbf{u}_d 的维数。概率密度的等值曲线可描述为:

$$(\xi - \mathbf{u}_d)^T R_{uu}^{-1} (\xi - \mathbf{u}_d) = \kappa \quad (9)$$

其中, κ 为任意正常数, 由可以确定曲线表面所包围的 n 维区域大小。本文中 $n=3$, 其表面为椭圆柱。

目标估计值 $\hat{\mathbf{u}}_d$ 位于椭圆柱内部的概率为:

$$\Pr(\kappa) = \iint_{\Omega} \cdots \int p_{\hat{u}_d}(\xi) \cdot d\xi_1 d\xi_2 \cdots d\xi_n \quad (10)$$

其中, 积分区域 $\Omega = \{\xi | (\xi - \mathbf{u}_d)^T R_{uu}^{-1} (\xi - \mathbf{u}_d) \leq \kappa\}$ 。为将多重积分转换为单重积分, 引入变量 $\gamma = \xi - \mathbf{u}_d$, 则可将式(10)转换为:

$$\Pr(\kappa) = \eta \cdot \iint_{\Omega_1} \cdots \int \exp\left\{-\frac{\gamma^T R_{uu}^{-1} \gamma}{2}\right\} d\gamma_1 d\gamma_2 \cdots d\gamma_n \quad (11)$$

其中, $\eta = (2\pi)^{-n/2} (\det [R_{uu}])^{-1/2}$, 积分区域 $\Omega_1 = \{\gamma | \gamma^T R_{uu}^{-1} \gamma \leq \kappa\}$ 。由于 R_{uu}^{-1} 是正定矩阵, 则一定存在正交矩阵 U 满足:

$$U^T R_{uu}^{-1} U = \Sigma^{-1} \Leftrightarrow R_{uu}^{-1} = U \Sigma^{-1} U^T \quad (12)$$

根据椭圆柱体积证明过程^[20], 得到:

$$\Pr(\kappa) = \text{erf}\left(\frac{\sqrt{\kappa/2}}{\sqrt{\pi}}\right) - \sqrt{2\kappa/\pi} \cdot \exp\left\{-\kappa/2\right\} \quad (13)$$

其中, $\text{erf}(\cdot)$ 为误差函数, 其定义为 $\text{erf}(x) = (2/\sqrt{\pi}) \int_0^x \exp(-t^2) \cdot dt$ 。

下面进行定位误差椭圆概率有效性分析, 假设 5 个测向站的经纬度坐标分别为 (116.47°E, 39.72°N)、(125.52°E, 44.24°N)、(113.55°E, 23.23°N)、(119.49°E, 39.38°N) 和 (121.53°E, 31.46°N), 目标位置 (134°E, 28°N)。测向站测向误差 σ_θ 服从零均值高斯分布, σ_θ 设为 0.3°、0.7° 和 1.0°, 采用交叉定位的方式求得目标估计位置, 将交点平均值 $\bar{\mathbf{u}}_d$ 替代目标真实值 \mathbf{u}_d , 进行 2000 次蒙特卡洛仿真实验, 则误差椭圆概率计算公式为:

$$\Pr(\kappa) = \Pr\left[(\hat{\mathbf{u}}_d - \bar{\mathbf{u}}_d)^T R_{uu}^{-1} (\hat{\mathbf{u}}_d - \bar{\mathbf{u}}_d)\right] \quad (14)$$

如图 4 所示, 误差椭圆概率随着参数 κ 的增大而增大, 仿真值和理论值吻合较好, 从而验证了理论分析的有效性。可以得出, 误差椭圆概率越大, 则定位精度越低; 误差椭圆概率越小, 定位精度越高。因此, t 时刻第 i 个智能体对信号定位的奖励为:

$$r_{t1}^i = 100 \times (1 - \Pr(\kappa)) \quad (15)$$

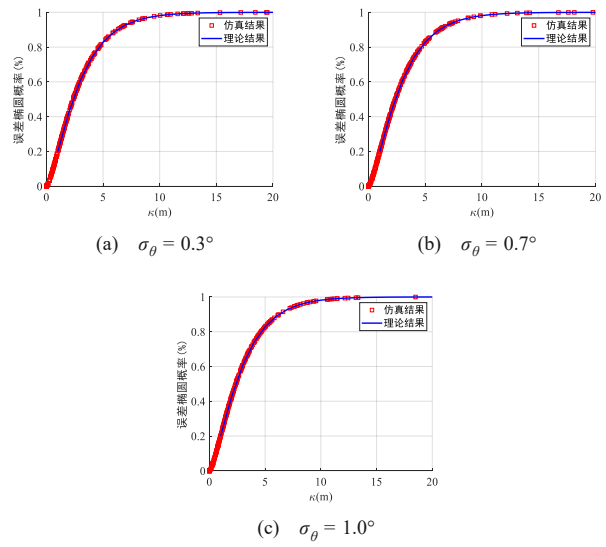


图 4 误差椭圆概率仿真结果与理论结果对比

(2) 检测奖励: 在信号检测任务中, 考虑 2GALE、110A、CLOVER2000、LINK4A 和 LINK11 的 5 类短波信号。需要说明的是, 本文构造的检测奖励模型是可根据信号类型的增多进行扩展的, 但是为了方便表述, 设定 5 类短波信号。如图 5 所示, 采用双塔网络计算信号类型匹配度, 双塔网络有两个相同的子网络组成, 两个网络共享权重, 确保提取的特征具有一致性。第一个网络的输入为智能体当前状态。第二个网络的输入为 5 种类型的短波信号时频图, 经过卷积层、池化层等特征提取模块, 得到特征向量 \mathbf{v}_1 和 \mathbf{v}_2^i 。然后, 两个向量计算余弦相似度作为匹配度, 取匹配度的最大值作为智能体 i 的检测奖励 r_{t2}^i 为:

$$r_{t2}^i = \max \frac{\mathbf{v}_1 \cdot \mathbf{v}_2^i}{|\mathbf{v}_1| \cdot |\mathbf{v}_2^i|} \quad (16)$$

(3) 协作奖励: 为避免智能体对信号的重复检测与定位, 如果 t 时刻信号 j 被标记为检测定位已完成, 则记为 $H_j(t) = 1$, 此时智能体 i 对信号 j 不再进行信号检测与定位; 否则 $H_j(t) = 0$ 。为实现多智能体协作, 智能体 i 在确认信号 j 未被其他智能体检测时, 可以获得额外协作奖励 r_{t3}^i 。

另外, 为保证智能体移动的有效性, 考虑对智能体 i 的每次移动增加移动惩罚 r_{t4}^i 。综上, t 时刻第 i 个智能体奖励为:

$$r_t^i = \sum_{i=1}^N \delta r_{t1}^i + (1 - \delta)r_{t2}^i + r_{t3}^i + r_{t4}^i \quad (17)$$

其中, $\delta \in [0,1]$ 为定位奖励权重。因此, 多智能体环境中的总奖励为 $\mathbf{R}_t = [r_t^1, r_t^2, \dots, r_t^n]$ 。

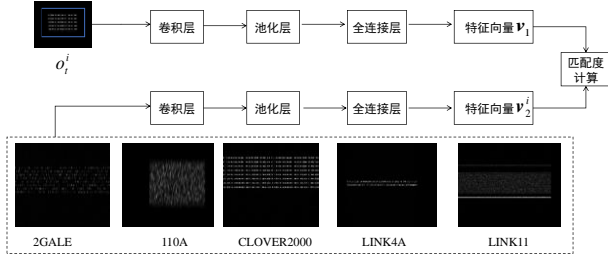


图5 信号匹配度计算

2.2 算法结构

MAPPO-DL 通过扩展 PPO 的框架, 结合集中训练与分布执行(Centralized Training with Decentralized Execution, CTDE)模式, 有效解决多智能体协作或竞争问题。如图 6 所示, MAPPO-DL 算法的实现涉及每个智能体训练两个网络, 即 Actor 网络和 Critic 网络。Actor 网络则根据贝尔曼方程得到动作价值函数 $Q(o_t^i, a_t^i)$, 学习一个从观测 o_t^i 到动作 a_t^i 的映射函数, 记为:

$$Q(o_t^i, a_t^i) = \mathbb{E} \left[r_t^i + \gamma \max_{a_t^i} Q(o_{t+1}^i, a_{t+1}^i) | o_t^i, a_t^i \right] \quad (18)$$

Critic 网络的目标是学习一个从观测空间 \mathbf{O}_t 到价值的映射函数 $V_\phi(\mathbf{O}_t)$, 记为:

$$V_\phi(\mathbf{O}_t) = \mathbb{E} \left[r_t^i + \gamma V_\phi(\mathbf{O}_{t+1}) | \mathbf{O}_t \right] \quad (19)$$

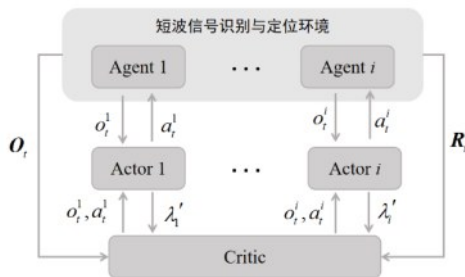


图6 MAPPO-DL 算法网络结构

在 MAPPO-DL 中, 每个智能体 i 的策略由 Actor 网络决定, Actor 网络更新基于以下目标函数:

$$L_i(\lambda_i) = \mathbb{E} \left[\min(b_t^i(\lambda_i) A_t^i, \text{clip}(b_t^i(\lambda_i), 1 - \varepsilon, 1 + \varepsilon) A_t^i) \right] \quad (20)$$

其中, $b_t^i(\lambda_i) = \pi_{\lambda_i}(a_t^i | o_t^i) / \pi_{\lambda_{i,\text{old}}}(a_t^i | o_t^i)$, ε 为裁剪率, $\text{clip}(x, l, r) = \max(\min(x, r), l)$, 即把 x 限制在 $[l, r]$ 内。在短波检测与定位这一特定场景下, 策略梯度更新幅度过大会导致学习过程不稳定。通过引入裁剪机制, 可以控制新旧参数之间的差异, 在信任域内搜索最优策略, 从而提升训练稳定性。 A_t^i 由集中式 Critic 网络计算。Critic 网络输入为全局状态 \mathbf{O}_t , 输出为全局优势值:

$$A_t^i = Q(o_t^i, a_t^i) - V_\phi(\mathbf{O}_t), a_t \sim \pi(a_t, o_t), \quad (21)$$

$$o_{t+1} \sim P(o_{t+1} | o_t, a_t)$$

其中, γ 为折扣回报率, $p(o_{t+1} | o_t, a_t)$ 为状态转移概率。另外, 为了防止智能体陷入次优策略, 在 Actor 网络的损失函数 $L_i(\lambda_i)$ 中加入了策略熵项, 并乘以系数 μ , 记为:

$$L'(\lambda_i) = L(\lambda_i) - \mu \sum_{a_t^i} \pi_{\lambda_i}(a_t^i | o_t^i) \log(\pi_{\lambda_i}(a_t^i | o_t^i)). \quad (22)$$

Critic 网络 V_ϕ 通过最小化均方误差更新目标函数:

$$L(\phi) = \mathbb{E} \left[(r_t^i + \gamma V_\phi(\mathbf{O}_{t+1}) - V_\phi(\mathbf{O}_t))^2 \right] \quad (23)$$

MAPPO-DL 使用卷积网络处理时频图像数据, 在时间维度上保留动态特征。如图 7 所示, MAPPO-DL 算法的 Actor 网络基于局部观测 $o_t^i(t)$ 生成动作 $a_t^i(t)$; Critic 网络利用全局状态 s_t 优化价值估计 V_ϕ , 如图 8 所示。执行时, 各智能体仅依赖局部观测独立决策。连续动作部分输出均值向量, 离散动作部分输出多个动作的概率分布。

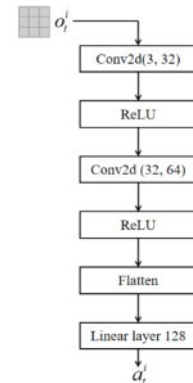


图7 Actor 网络结构

2.3 算法流程

算法 1 展示了 MAPPO-DL 算法流程。具体来说, MAPPO-DL 算法首先初始化 Actor 网络和

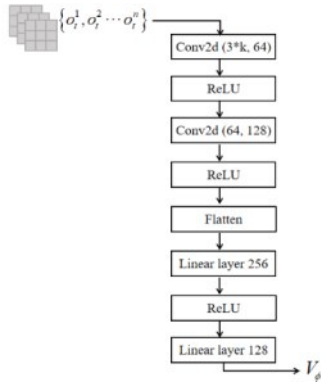


图 8 Critic 网络结构

Critic 网络的参数，以减轻训练初期出现的梯度消失或爆炸等问题。

在训练阶段， N 个智能体每回合迭代收集长度为 T 的步长。智能体 i 选中测向时频图数据获取局部观测 o_t^i ，依据 Actor 网络选择动作 a_t^i ，实现智能体与环境的交互。其次，每个智能体在交互过程中获取并存储观测状态、动作、当前奖励、下一状态，记作 $(\mathbf{O}_t, \mathbf{A}_t, \mathbf{R}_t, \mathbf{O}_{t+1})$ 。

基于收集到的经验回放数据集，采用广义优势估计方法来估算优势函数 A_t^i ，利用值归一化方法来计算归一化的值函数 V_ϕ 。最后，Actor 网络与 Critic 网络均采用 Adam 优化器进行参数更新，并配合梯度裁剪与衰减学习率，其中梯度更新的幅度由参数 τ 限定，以避免梯度爆炸问题。训练过程持续进行，直到完成预定的训练轮次 E 。

2.4 算法性能分析

在 MAPPO 算法中，采用裁剪操作来限制策略梯度更新的幅度。通过使用裁剪操作，MAPPO 与传统的信赖域策略优化(Trust Region Policy Optimization, TRPO)中使用的 KL (Kullback-Leibler) 散度约束有所不同，从而显著提高了算法的灵活性^[21-22]。设 π_{λ} 表示一个随机策略，MAPPO 通

过裁剪 $A_{\pi_{\lambda}}$ 来保证 $\pi_{\lambda_{old}}$ 和 $\pi_{\lambda_{new}}$ 之间的差异。因此，MAPPO-DL 的目标函数为：

$$J(\lambda) = \mathbb{E}_{\pi_{\lambda_{old}}} \left[\min \left\{ \frac{\pi_{\lambda_{new}}}{\pi_{\lambda_{old}}} A_{\lambda_{old}}, \text{clip}_{1-\varepsilon}^{1+\varepsilon} \left(\frac{\pi_{\lambda_{new}}}{\pi_{\lambda_{old}}} \right) A_{\lambda_{old}} \right\} \right] \quad (24)$$

1) 收敛性分析

验证 MAPPO-DL 的学习收敛性能，即验证当满足裁剪操作中的分布约束：

算法 1: MAPPO-DL

输入：短波信号检测与定位时频图环境

输出：信号估计类型 Z 和估计位置 $\hat{\mathbf{u}}_d$

1) 初始化经验回放池 B ，初始化 Critic 网络 $V_\phi(a_1, a_2, \dots, a_N)$ 的参数 ϕ ，初始化每个智能体 Actor 网络 $\pi_{\lambda_i}(o_t^i)$ 的参数 λ_i ，初始化目标 Critic 网络 $V_\phi'(a_1, a_2, \dots, a_N)$ 的参数 ϕ' ，设置裁剪率 ε ，Actor 网络学习率 α ，Critic 网络学习率 β 。

2) for $EP=1:E$ do

3) 滤波窗智能体和目标信号随机初始位置，获得初始观测状态 s_0 和各智能体观测 $\{o_0^1, o_0^2, \dots, o_0^N\}$ 。

4) for $SP=1:T$ do

5) 每个智能体根据局部观测状态 o_t^i ，依据 Actor 网络 $\pi_{\lambda_i}(o_t^i)$ 选择动作 a_t^i 。

6) 执行动作 $a_t = \{a_t^1, a_t^2, \dots, a_t^N\}$ 。

7) 依据选中时频图 MUSIC 测向结果 θ_i 计算定位奖励 r_{t1}^i ，并得到信号估计位置 $\hat{\mathbf{u}}_d$ 。

8) 根据信号类型检测结果 Z 计算 r_{t2}^i ，得到总奖励 r_t^i ，并获得新的智能体观测 o_{t+1}^i 。

9) 计算广义优势估计 $A_t^i = R_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t)$ 。

10) 将 $(\mathbf{O}_t, \mathbf{A}_t, \mathbf{R}_t, \mathbf{O}_{t+1})$ 存入经验回放池 B 。

11) 计算 Critic 网络损失 $L(\phi) = \mathbb{E} \left[(R_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t))^2 \right]$ 。

12) 计算 Actor 网络损失 $L_i(\lambda_i) = \mathbb{E} \left[\min(b_{i,t}(\lambda_i) A_{i,t}, \text{clip}(b_{i,t}(\lambda_i), 1 - \varepsilon, 1 + \varepsilon) A_{i,t}) \right]$ 。

13) 更新 Actor 网络参数 $\lambda_i \leftarrow \text{Adam}(\lambda_i, \nabla_{\lambda_i} L_i(\lambda_i))$ 。

14) 更新 Critic 网络参数 $\phi \leftarrow \text{Adam}(\phi, \nabla_{\phi} L(\phi))$ 。

15) 更新目标 Critic 网络 $\phi' \leftarrow \tau \times \phi + (1 - \tau) \times \phi'$ 。

16) end for

17) end for

$$1 - \varepsilon \leq \frac{\pi_{\lambda_{new}}}{\pi_{\lambda_{old}}} \leq 1 + \varepsilon \quad (25)$$

时，MAPPO 中策略更新的结果与初始策略梯度方法的结果之间存在约束变化。设 P_π 为折扣访问频率：

$$P_\pi(s) = P(s_0 = s) + \gamma P(s_1 = s) + \gamma^2 P(s_2 = s) + \dots \quad (26)$$

其中， $P(s_i = s)$ 表示在采样轨迹的第 i 个时刻遇到状态 s 的概率。为了表示新的策略 $\pi_{\lambda_{new}}$ 预期价值收益相对旧的策略 $\pi_{\lambda_{old}}$ 优势，则有：

$$\begin{aligned}
V(\pi_{\lambda_{new}}) &= V(\pi_{\lambda_{old}}) + \mathbb{E}_{s_0, a_0, \dots, \pi_{\lambda_{new}}} \left[\sum_{t=0}^{\infty} \gamma^t A_{\pi_{\lambda_{old}}}(s_t, a_t) \right] \\
&= V(\pi_{\lambda_{old}}) + \sum_s P_{\pi_{\lambda_{new}}}(s) \sum_a \pi_{\lambda_{new}}(a|s) A_{\pi_{\lambda_{old}}}(s, a)
\end{aligned} \quad (27)$$

上式表明, 如果策略更新 $\pi_{\lambda_{old}} \rightarrow \pi_{\lambda_{new}}$ 在每个状态 s 下都能带来正向的价值期望, 即 $\sum_a \pi_{\lambda_{new}}(a|s) A_{\pi_{\lambda_{old}}}(s, a) \geq 0$, 则策略价值函数 V_ϕ 就不会下降。然而, 在近似过程中, 由于估计和近似误差, 通常不可避免地会存在一些状态 s , 对于这些状态的价值期望为负, 即 $\sum_a \pi_{\lambda_{new}}(a|s) A_{\pi_{\lambda_{old}}}(s, a) < 0$ 。由于 $P_{\pi_{\lambda}}(s)$ 对 $\pi_{\lambda_{new}}$ 的相互依赖, 直接优化公式(24)比较困难。因此, 考虑引入局部近似方法:

$$\begin{aligned}
G_{\pi_{\lambda_{old}}}(\pi_{\lambda_{new}}) &= V(\pi_{\lambda_{old}}) + \\
&\sum_s P_{\pi_{\lambda_{old}}}(s) \sum_a \pi_{\lambda_{new}}(a|s) A_{\pi_{\lambda_{old}}}(s, a)
\end{aligned} \quad (28)$$

其中使用 $\pi_{\lambda_{old}}$ 代替 $\pi_{\lambda_{new}}$ 。在此基础上, 定理 1 给出了策略价值函数 V_ϕ 改进的下界。

定理 1 通过调整 MAPPO 策略更新过程中 $\pi_{\lambda_{new}}$ 和 $\pi_{\lambda_{old}}$ 之间的差异, 改进的 V_ϕ 下界可以表示为:

$$G_{\pi_{\lambda_{old}}}(\pi_{\lambda_{new}}) - V(\pi_{\lambda_{old}}) \geq \frac{4\gamma\rho}{(1-\gamma)^2} \alpha^2 \quad (29)$$

其中, α 表示分歧的概率, $\rho = \max |A_{\pi_{\lambda}}(s, a)|$, 具体证明见附录 I。

2) 计算复杂度分析

设环境中存在 N 个智能体滤波窗, 每个智能体的 Actor 和 Critic 网络均为 CNN, 网络层数为 L , 卷积核大小固定, 每层输出通道数为 C_l , 输入时频图

尺寸为 $H \times W$ 。CNN 的浮点运算次数约为 $O(\sum_{l=1}^L C_{l-1} \cdot k^2 \cdot H_l \cdot W_l)$, 其中 k 为卷积核大小, H_l, W_l 为第 l 层特征图尺寸。由于采用全卷积结构, 计算量与输入图像尺寸近似成线性关系, 记为 $O(\mathcal{F}(H, W))$ 。

每个智能体需完成一次 Actor 前向、Critic 前向、优势估计、网络反向传播。其中反向传播计算量约为前向的 2~3 倍, 因此单智能体单步训练时间复杂度为 $O(\mathcal{F}(H, W))$ 。

由于 MAPPO-DL 中多智能体共享网络参数, N 个智能体在同一环境步内的计算可并行, 故单步总时间复杂度仍为 $O(\mathcal{F}(H, W))$, 与 N 无关。训练需经历 E 个回合, 每回合最大步长为 T , 总时间复杂度为 $O(E \cdot T \cdot \mathcal{F}(H, W))$ 。

3 仿真实验及结果分析

3.1 仿真环境设置

本文实验通过模拟产生 5 类短波信号, 并添加噪声, 然后进行 FFT 处理得到时频图。在信号产生阶段, 设定时频图持续时间为 4s, 信号持续时间为 100ms 至 500ms, 符号速率 f_b 为 2400 Bd - 4800Bd, 采样频率 f_s 为 120kHz, FFT 点数 N_{FFT} 为 1024, 载波频率 f_c 为 25 kHz。为抑制码间串扰并提高频谱利用率, 采用平方根升余弦滚降滤波器进行脉冲成型, 余弦滚降系数 α 为 0.8, 信号带宽 B 为 4.32 kHz - 8.64 kHz, SNR 范围为 5 dB - 25dB。

强化学习环境为随机信噪比的时频图, 在上面添加 2 至 6 个随机分布的短波信号。每张时频图中短波信号的类型、信噪比、带宽、幅度和符号速率等参数均在特定范围内随机变化。时频图大小固定

表 1 信号和强化学习环境仿真参数设置

| 符号 | 描述 | 数值 | 符号 | 描述 | 数值 |
|----------|--------|---------------------|------------|----------|---------------|
| f_b | 信号符号速率 | 2400 Bd - 4800Bd | t_s | 最大信号持续时间 | 100ms - 500ms |
| f_s | 采样频率 | 25 kHz | w_{\max} | 智能体的最大宽度 | 100 |
| α | 余弦滚降系数 | 0.8 | w_{\min} | 智能体的最小宽度 | 10 |
| B | 信号带宽 | 4.32 kHz - 8.64 kHz | h_{\max} | 智能体的最大高度 | 100 |
| SNR | 信号信噪比 | 5 dB - 25dB | h_{\min} | 智能体的最小高度 | 5 |

为 400×300 , 利用 `pygame` 创建时频图强化学习环境, 并随机产生矩形滤波窗表示智能体。信号参数和强化学习环境设置如表 1 所示。

站点与目标均在 $(10^\circ\text{N}, 100^\circ\text{E}) \times (40^\circ\text{N}, 140^\circ\text{E})$ 的地理区域内选取。环境依赖关系根据第 3 节所述进行 MDP 建模。本实验所用 `torch` 版本为 1.13.0, 计算机显卡配置为 NVIDIA GeForce RTX A5000。训练时, 每一回合随机初始智能体和信号位置, 时频图上的不同信号代表不同的目标。

另外, 需要对奖励函数中的超参数 δ 进行取值测试, 以保证所提算法的最佳检测与定位性能。图 9 展示了奖励函数值随 δ 的变化情况。可以看出, 在 δ 值为 0.7 时, 定位奖励、检测奖励、协作奖励和综合奖励值最大。因此, 在后续实验中, 为获得更好的训练效果, 考虑 δ 取值为 0.7。

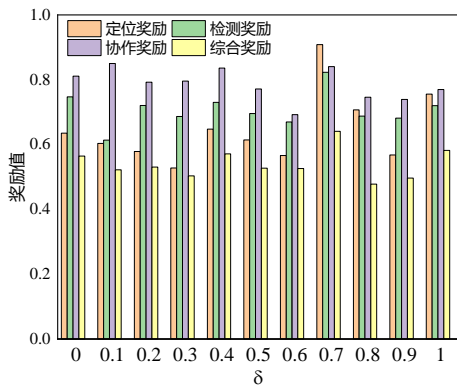


图 9 各奖励值随 δ 变化情况

3.2 基线算法

为验证本文 MAPPO-DL 算法与所设计奖励函数的有效性, 本文分别采用 4 种 MARL 基线算法进行对比, 分别是多智能体深度确定性策略梯度算法 (Multi-Agent Deep Deterministic Policy Gradient, MADDPG)^[23]、基于注意机制的多智能体演员评论家算法 (Multi-Actor-Attention-Critic, MAAC)^[24]、查询-策略对齐算法 (Query-Policy Alignment, QPA)^[25] 和近端策略探索 (Proximal Policy Exploration, PPE)^[26]。

另外, 为进一步验证所提 MAPPO-DL 算法对短波信号检测与定位的有效性, 考虑设计 3 种算法对检测与定位模块进行测试, 分别是:

1) 基于目标检测的多智能体深度强化学习 (MAPPO-D): 奖励函数中只考虑目标检测结果,

即定位奖励权重 $\delta = 0$ 。

2) 基于目标定位的多智能体深度强化学习 (MAPPO-L): 奖励函数中只考虑目标信号定位结果, 即定位奖励权重 $\delta = 1$ 。

3) 基于目标检测与定位的单智能体深度强化学习 (PPO-DL): 算法中只有一个智能体, 奖励函数中同时考虑目标信号检测和定位情况。

所提算法与基线算法训练过程的参数设置如表 2 所示。

表 2 强化学习训练参数设置

| 符号 | 描述 | 数值 |
|------------|---------|--------|
| γ | 折扣回报率 | 0.95 |
| ϵ | 裁剪率 | 0.2 |
| α | 价值网络学习率 | 0.0001 |
| β | 策略网络学习率 | 0.0003 |
| B | 经验回放池大小 | 1000 |
| E | 回合数 | 5000 |
| T | 执行最大步长 | 50 |
| r_{i3}^i | 协作奖励 | 20 |
| r_{i4}^i | 移动惩罚 | -10 |

3.3 训练结果对比

图 10 展示了 MAPPO-DL 算法和不同 MARL 算法的平均回报收敛性曲线和检测定位时间的对比曲线, 阴影部分表示训练结果标准差。

图 10 (a) 可以看出 MAPPO-DL 算法收敛后的平均回报最高, 且收敛较快。图 10 (b) 展示了 MAPPO-DL 算法收敛后的平均检测定位时间最低, 且训练过程更稳定。而 MADDPG 和 MAAC 都是离线更新算法, 在多智能体信号检测与定位环境训练中收敛较慢, 收敛后的平均奖励也低于所提算法。QPA 算法更专注于在当前策略分布内学习奖励, 但会依赖于人类反馈。PPE 算法虽然可扩展在缺乏采样的策略近端区域的探索范围, 但它对经验缓冲区的要求更高。

图 11 展示了 MAPPO-DL 算法和 MAPPO-D、MAPPO-L 和 PPO-DL 的平均回报收敛性曲线和检测定位时间的对比曲线。

图 11 (a) 可以看出四种强化学习算法均能在多次训练后收敛。从收敛后平均回报来看, MAPPO-DL 算法收敛后的平均回报最大, 其次是 PPO-DL

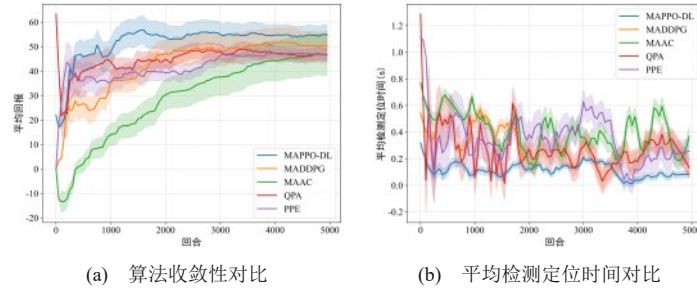


图 10 不同MARL算法训练结果对比

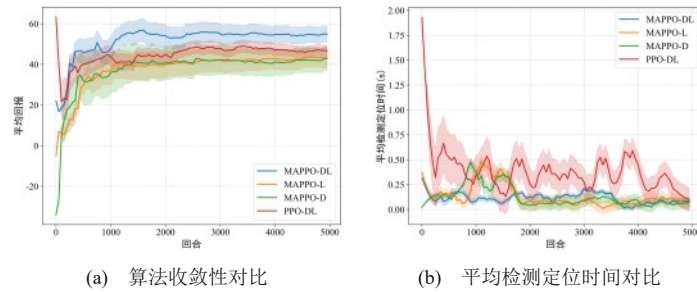


图 11 所提算法不同模块训练结果对比

算法，原因在于PPO-DL算法的奖励函数也考虑了目标信号检测和定位情况；MAPPO-D和MAPPO-L收敛后的平均回报基本一致。从收敛速度来看，四种算法收敛速度相当，但仍可以比较出MAPPO-D和MAPPO-L收敛较快，MAPPO-DL收敛较慢。原因在于多智能体之间的协作影响了收敛速度。同时裁剪操作保证了策略更新，使得奖励曲线较为平滑，避免了震荡。

图 11 (b)可以看出多智能体算法检测定位时间均低于单智能体算法，MAPPO算法的波动较小。而PPO-DL算法的检测定位时间较长且波动较大，原因在于单智能体强化学习在处理多目标问题时，需要在时频图上逐一寻找信号，再完成检测与定位工作。

由于多智能体强化学习方法中的协作策略可以有效共享智能体之间的信息，进一步协调智能体的动作。因此，MAPPO-DL克服了单智能体深度强化学习算法的局限。

4 模型性能测试

在完成深度强化学习算法训练后，选择MAPPO-DL和基线算法第5000步训练得到的策略网络参数作为测试初始模型，在不同的多目标场景进行500次信号检测与定位，对比分析目标的定位

时间、目标检测率和定位误差。

4.1 关于信号检测与定位时间的性能测试

图 12 表示算法测试时间和测试平均奖励的变化。考虑在测试集上运行500次，记录从输入时频图到输出检测定位结果的平均时间。如表3所示，MAPPO-DL因采用并行计算与参数共享，单步训练时间低于MADDPG、MAAC、QPA、PPE独立更新算法；收敛总时间虽略高于MAPPO-D和MAPPO-L，但测试时间明显低于其他算法。

表 3 不同算法信号检测与定位时间

| 算法 | 单步训练时间(ms) | 测试时间(ms) |
|----------|------------|----------|
| MAPPO-DL | 12.1 | 0.02 |
| MAPPO-L | 10.2 | 0.12 |
| MAPPO-D | 10.3 | 0.10 |
| PPO-DL | 9.7 | 0.14 |
| MADDPG | 17.4 | 0.08 |
| MAPPO | 16.2 | 0.06 |
| QPA | 17.8 | 0.08 |
| PPE | 19.6 | 0.09 |

从测试时间可以看出，图 12 (c)中MAPPO-DL算法的测试时间最短，平均测试时间为0.02s；图

12(d)中 PPO-DL 算法的测试时间最长;图 12 (a)中 MAPPO-L 算法与图 12 (b)中 MAPPO-D 算法定位时间相当,平均测试时间为 0.13s。

同样,单智能体的 PPO-DL 算法在处理多信号检测与定位问题时,需要逐一测试,导致测试信号检测与定位的时间较长,平均测试时间为 0.15s。因此,MAPPO-DL 平均测试时间约 0.02s,MAPPO-DL 算法将信号检测与定位时间平均缩短了 0.12s。MAPPO-DL 算法较其他多智能体算法更优,满足短波信号实时检测需求。

从测试平均奖励来看,四种算法的测试平均奖励均保持相对稳定,这说明保存的算法模型已得到稳定的训练。相比于单智能体强化学习方法,多智能体强化学习方法无需逐一检测信号,通过智能体之间的协作完成多目标的实时检测与定位工作。

4.2 关于定位误差的性能测试

在进行定位误差的测试时,考虑使用的测向站

点与目标经纬度坐标如表 4 所示。测向站和目标的相对位置如图 13 所示。假设使用三个场景测试定位误差,分别是:

- (1) 场景 1: 5 个测向站和 2 个目标,测向站为 1、2、3、4 和 5,目标选择目标 1 和目标 2。
- (2) 场景 2: 5 个测向站和 2 个目标,测向站为 1、2、3、4 和 5,目标选择目标 3 和目标 4。
- (3) 场景 3: 5 个测向站和 3 个目标,测向站为 1、2、3、4 和 5,目标选择目标 5、目标 6 和目标 7。

为进一步突出多智能体的协作优势,本节仅使用 MAPPO-DL 网络模型在 3 种场景下进行定位性能测试。由于目标定位结果是依据各个测向站交点平均值作为目标估计值,因此目标定位结果依赖于 MUSIC 算法的测向精度。进一步地得出,目标定位精度与信号信噪比 SNR 成正比。因此,本节考虑采用均方根误差(Root mean square error,

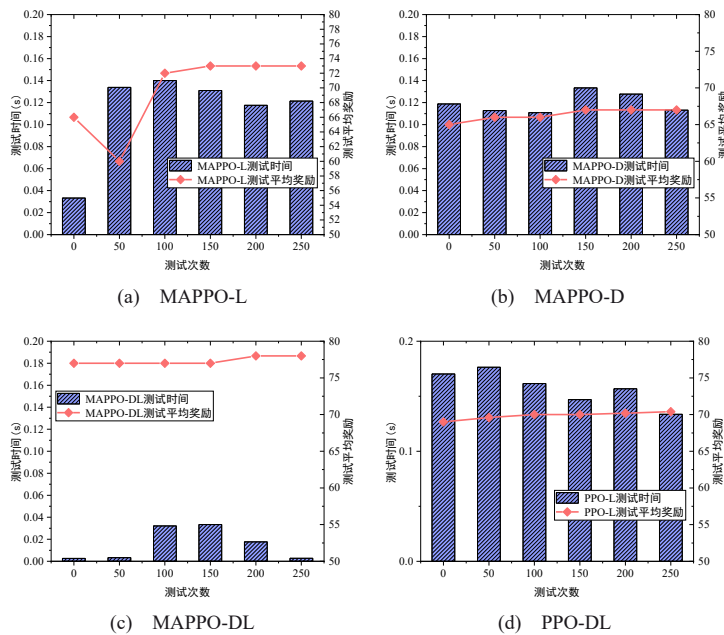


图 12 不同算法测试时间和测试平均奖励

表 4 测向站点与目标经纬度坐标

| 类型 | 站点 1 | 站点 2 | 站点 3 | 站点 4 | 站点 5 | 目标 1 |
|-------|---------|---------|---------|---------|---------|---------|
| 经度(E) | 117.00° | 106.82° | 104.95° | 118.18° | 124.17° | 103.40° |
| 纬度(N) | 36.30° | 14.75° | 43.10° | 25.54° | 43.05° | 34.69° |
| 类型 | 目标 2 | 目标 3 | 目标 4 | 目标 5 | 目标 6 | 目标 7 |
| 经度(E) | 97.40° | 101.40° | 95.70° | 98.57° | 92.82° | 96.50° |
| 纬度(N) | 35.42° | 32.69° | 33.32° | 30.26° | 31.54° | 28.45° |

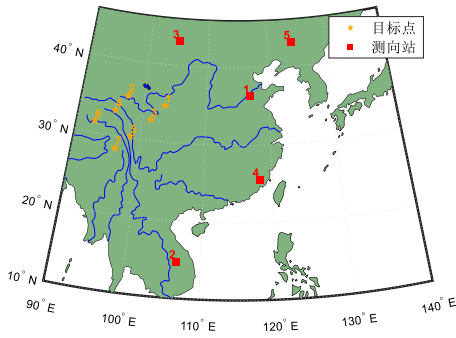


图 13 测向站与目标相对位置

RMSE)来对比算法对目标的定位性能，计算方式为：

$$RMSE(\mathbf{u}) = \sqrt{\frac{1}{K} \sum_{k=1}^K \|\hat{\mathbf{u}}_k - \mathbf{u}_k\|_2^2} \quad (30)$$

其中， $K=2000$ 为蒙特卡洛次数， $\hat{\mathbf{u}}_k$ 为第 k 次测试中目标位置的估计值， \mathbf{u}_k 为目标位置的真实值。

图 14 展示了三种场景下，不同目标随信噪比的变化，可以看出随着信噪比的增加，定位误差逐渐减小。如图 14 (a)场景 1 中，低信噪比时目标 1 的定位误差大于目标 2，高信噪比时二者的定位误差相当。如图 14 (b)的场景 2，低信噪比时目标 3 的定位误差小于目标 4。如图 14 (c)的场景 3，低信噪比时目标 6 的定位误差最大，目标 5 的定位误差最小，高信噪比时三者的定位误差相当。

因此，MAPPO-DL 算法能适应动态信号环境，具有较好的稳健性。通过以上分析可以得出，奖励函数中定位误差椭圆概率的设计可以有效引导智能体学习较低定位精度的最优策略。

进一步分析目标与测向站的平均距离对 MAPPO-DL 算法的影响。图 15 展示了 7 个目标与 5 个测向站之间的平均距离，目标 1、3、5 和 7 与测向站之间的平均距离较大，目标 2、4 和 6 与测向站之间的平均距离较小。另外，图 15 还展示了 7 个目标的 RMSE 的对比，可以看出，信噪比为 -10dB 时，目标 2、4 和 6 的 RMSE 明显大于目标 1、3、5、7。信噪比为 10dB 时，目标 2、4 和 6 的 RMSE 与目标 1、3、5、7 相差不大；随着目标与测向站的平均距离增大，RMSE 逐渐增大，且低信噪比时的变化趋势较高信噪比明显。因此，低信噪比时，MAPPO-DL 的定位精度对目标与测向站的远近较敏感。

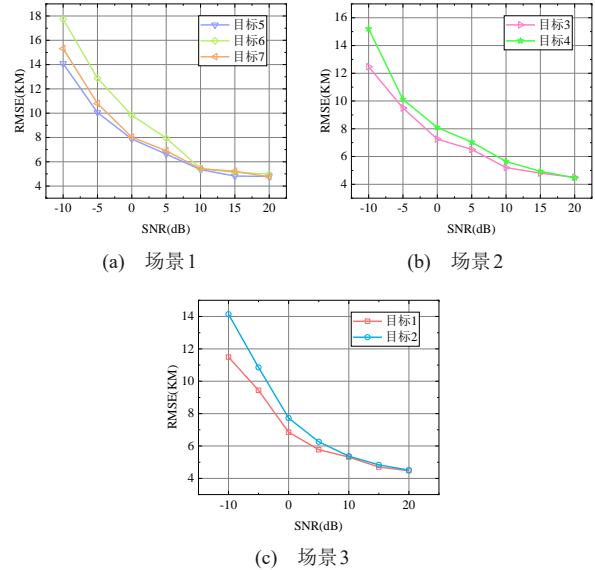


图 14 不同场景定位误差随信噪比变化

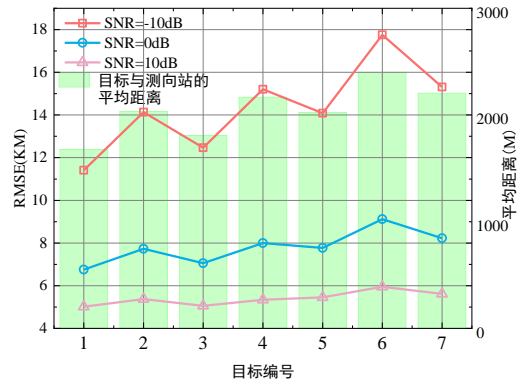
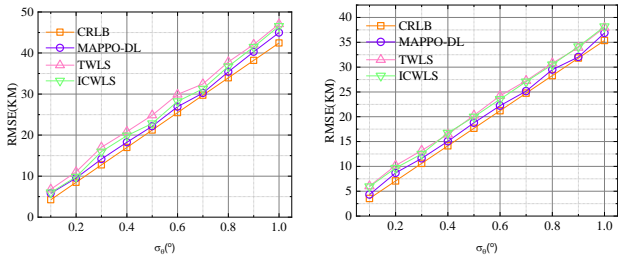


图 15 目标在不同信噪比下定位误差变化

假设目标 1、3、5、7 的位置为某一移动目标 A 的轨迹，目标 2、4 和 6 为某一移动目标 B 的轨迹，图 16(a)和图 16(b)分别展示了目标 A 的轨迹和目标 B 的轨迹在不同测向误差下所提算法与 TWLS 算法^[1]、ICWLS 算法^[27]的 RMSE 以及克拉美罗下界(Cramer-Rao lower bound, CRLB)的变化。可以看出，所提 MAPPO-DL 算法的 RMSE 与其他两种传统方法相比更小，与 CRLB 更接近。通过上述实验，验证了所提算法在信号时变条件下的检测与定位的鲁棒性。

4.3 关于多目标信号检测与定位准确率的性能测试

多目标信号检测与定位的准确率可以反映算法的检测正确率和定位精度。目标信号检测与定位的准确率 P_i 定义为：



(a) 目标 A 的定位误差随 σ_θ 变化 (b) 目标 B 的定位误差随 σ_θ 变化
图 16 目标 A 和目标 B 的轨迹在不同测向误差 σ_θ 下定位误差变化

$$P_l = \frac{N_{DL}}{N} \quad (31)$$

其中, $N = 500$ 表示测试样本数, N_{DL} 表示信号检测与定位成功的样本数。

为全面评估 MAPPO-DL 算法在不同信噪比条件下的定位精度, 固定 $\sigma_\theta = 0.5^\circ$ 采用平均精度均值 (mean Average Precision, mAP)^[28] 评价信号检测与定位精度。即定义多个定位误差阈值计算各定位误差阈值下的定位准确率, 并取其平均值作为综合定位精度指标。

如图 17 所示, 当定位误差阈值大于 25km 时, mAP 范围为 80% 以上。因此, 设置当目标定位误差小于 25km 时判定为定位成功。需要说明的是, 短波定位的相对误差通常小于 3% 时, 即满足短波定位需求^[29-30]。本文实验场景中测向站与目标之间的平均距离为 2015.08km, 所采用阈值为 25km 误差在该距离上相对误差为 1.24%, 满足短波定位需求。

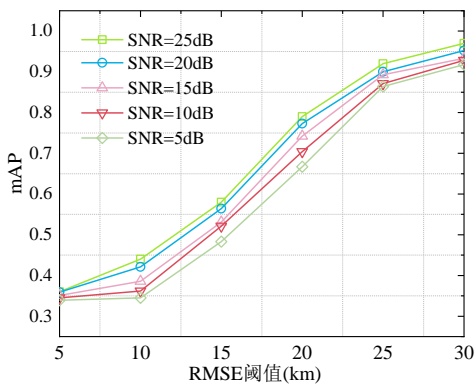
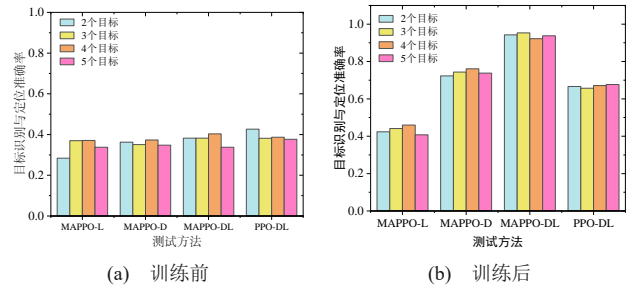


图 17 不同目标数量下的信号检测与定位准确率

图 18 表示训练前后模型对不同目标的信号检测率对比。图 18(a)中, 训练前四种算法在不同目标数量下的检测率均较低, 说明深度强化学习模型在与环境交互之前未学到策略, 只能随机的输出结

果。图 18(b)展示了训练后模型对不同目标的信号检测与定位准确率, MAPPO-DL 的目标检测与定位准确率最高, 达到 92%; 其次是 MAPPO-D 的目标检测与定位正确率为 70%, PPO-DL 的目标检测与定位正确率为 62%, MAPPO-L 的目标检测与定位准确率仅为 40%, 主要原因在于 MAPPO-L 算法在训练时奖励未考虑目标检测结果的约束, 因此测试时对信号检测的表现较差。



(a) 训练前 (b) 训练后
图 18 训练前后不同目标个数下的信号检测与定位准确率

相比目标检测与定位准确率最佳的基线算法, MAPPO-DL 算法的目标检测与定位准确率提升了 22%。信号检测与定位准确率的提高有效验证了信号匹配度奖励部分通过强化特征提取网络实现正确信号的检测。

4.4 关于多目标信号检测与定位可扩展性分析

本节旨在探究算法在更大规模协作场景下的扩展性。分别设置不同智能体数量, 固定不同目标数量, 训练至收敛, 并测试模型的信号识别与定位准确率。如图 19 所示, 由于采用集中训练分布执行, 随着目标数量的增加, 目标识别与定位准确率逐渐下降。所提 MAPPO-DL 目标识别与定位准确率明显高于其他 MARL 算法, MAPPO-DL 的策略熵可防止训练进入次优策略, 增加了算法的可扩展性。然而, 当目标数量超过 8 个时, MADRL 算法的性能会因维度灾难而显著下降, 主要原因为 Critic 网络直接拼接所有智能体观测的方式在高维输入下难以有效区分各智能体的贡献。

为缓解上述问题, 考虑在 MAPPO-DL 的 Critic 网络中引入多头注意力机制, 称为 MAPPO-DLA (MAPPO-DL with attention mechanism)。通过自适应学习不同智能体观测的权重, 更有效地提取协作信息, 缓解维度灾难带来的性能瓶颈。而 MAPPO-DLA 在 10 个目标的场景下, 比 MAPPO-DL 准确率

提升了约4%，可以有效抑制冗余信息干扰。因此，在基于强化学习的短波定位工作中引入注意力机制可以作为未来工作的重要方向之一。

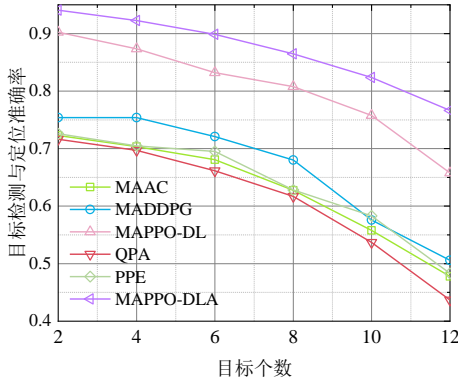


图 19 不同目标数量下的信号检测与定位准确率

5 结束语

本文提出了一种基于MAPPO的短波信号自主检测与定位方法，通过设计多智能体强化学习环境，运用深度强化学习算法探索试错，模拟短波信号自主检测与定位工作。MAPPO-DL可以快速实现目标的测向定位，MAPPO-DL在保持高检测定位精度的同时，具备较低的计算开销与良好的可扩展性，对后续精准定位目标具有一定参考价值。最后，在仿真阶段测试了MAPPO-DL方案的实时性和有效性，并且基于深度强化学习的方法不依赖大量人工标注，可以在线多回合自主学习。在未来工

作中，将继续研究基于深度强化学习的目标跟踪与轨迹分选任务。

附录I. 定理1证明

设 $D = \left(1 + \gamma P_{\pi_{\lambda_{old}}} + (\gamma P_{\pi_{\lambda_{old}}})^2 + \dots\right) = \left(1 - \gamma P_{\pi_{\lambda_{old}}}\right)^{-1}$, $\tilde{D} = \left(1 + \gamma P_{\pi_{\lambda_{new}}} + (\gamma P_{\pi_{\lambda_{new}}})^2 + \dots\right) = \left(1 - \gamma P_{\pi_{\lambda_{new}}}\right)^{-1}$ 。约定状态空间的密度 χ 是一个向量, 状态空间上的奖励函数 r 是一个对偶向量, 即向量上的线性泛函, 因此 $r\chi$ 是一个标量, 表示在密度 χ 下的期望奖励, 则 $V(\pi_{\lambda_{old}}) = rD\chi_0$, $V(\pi_{\lambda_{new}}) = c\tilde{D}\chi_0$ 。设 $\Delta = P_{\pi_{\lambda_{new}}} - P_{\pi_{\lambda_{old}}}$, $V(\pi_{\lambda_{new}}) - V(\pi_{\lambda_{old}}) = r(\tilde{D} - D)\chi_0$ 进行界定, 需要从一些标准的扰动理论进行推导:

$$D^{-1} - \tilde{D}^{-1} = \left(1 - \gamma P_{\pi_{\lambda_{old}}}\right) - \left(1 - \gamma P_{\pi_{\lambda_{new}}}\right) = \gamma \Delta \quad (32)$$

左乘 D , 右乘 \tilde{D} , 得到:

$$\tilde{D} - D = \gamma D \Delta \tilde{D} \Rightarrow \tilde{D} = D + \gamma D \Delta \tilde{D} \quad (33)$$

将右侧代入 \tilde{D} 得到:

$$\tilde{D} = D + \gamma D \Delta D + \gamma^2 D \Delta D \Delta \tilde{D} \quad (34)$$

可以得到:

$$\begin{aligned} V(\pi_{\lambda_{new}}) - V(\pi_{\lambda_{old}}) &= r(\tilde{D} - D)\chi \\ &= \gamma r D \Delta D \chi_0 + \gamma^2 r D \Delta D \tilde{D} \chi_0 \end{aligned} \quad (35)$$

下面首先考虑主导项 $\gamma r D \Delta D \chi_0$ 。显然 $rD = v$, 即无穷远状态值函数, 同时 $D\chi_0 = \chi_\pi$ 。因此, 可以得出 $\gamma v D \Delta D \chi_0 = \gamma v \Delta \chi_\pi$ 。下面将证明该表达式等于期望的优势 $L_\pi(\pi_{\lambda_{new}}) - L_\pi(\pi_{\lambda_{old}})$, 具体为:

$$\begin{aligned} L_\pi(\pi_{\lambda_{new}}) - L_\pi(\pi_{\lambda_{old}}) &= \sum_s \chi_\pi(s) \sum_a (\pi_{\lambda_{new}}(a|s) - \pi_{\lambda_{old}}(a|s)) A_\pi(s, a) \\ &= \sum_s \chi_\pi(s) \sum_a (\pi_\theta(a|s) - \pi_\theta(a|s)) \cdot \\ &\quad \left[r(s) + \sum_{s'} p(s'|s, a) \gamma v(s') - v(s) \right] \\ &= \sum_s \chi_\pi(s) \sum_{s'} \sum_a (\pi_{\lambda_{old}}(a|s) - \pi_{\lambda_{new}}(a|s)) \\ &\quad p(s'|s, a) \gamma v(s') \\ &= \sum_s \chi_\pi(s) \sum_{s'} (p_\pi(s'|s) - p_{\pi_{\lambda_{new}}}(s'|s)) \gamma v(s') \\ &= \gamma v \Delta \chi_\pi \end{aligned} \quad (36)$$

然后, 对 $O(\Delta^2)$ 项 $\gamma^2 r D \Delta D \Delta \pi_{\lambda_{new}} \chi$ 进行界定。首先, 考虑 $\gamma r D \Delta = \gamma v \Delta$ 和对偶向量的分量 s , 则有:

$$\begin{aligned}
\left| (\gamma v \Delta)_s \right| &= \left| \sum_a (\pi_{\lambda_{new}}(s, a) - \pi_{\lambda_{old}}(s, a)) Q_\pi(s, a) \right| \\
&= \left| \sum_a (\pi_{\lambda_{new}}(s, a) - \pi_{\lambda_{old}}(s, a)) A_\pi(s, a) \right| \quad (37) \\
&\leq \sum_a \left| \pi_{\lambda_{new}}(s, a) - \pi_{\lambda_{old}}(s, a) \right| \cdot \max_a |A_\pi(s, a)| \\
&\leq 2\alpha\rho
\end{aligned}$$

其中, 最后一步使用了总变差散度的定义, 以及 $\rho = \max |A_\pi(s, a)|$ 的定义, α 为不同意概率。通过 ℓ_1 算子范数来对另一部分 $D\Delta\tilde{D}\chi$ 进行有界处理, 记为:

$$\|A\|_1 = \sup_{\rho} \{ \} \quad (38)$$

由于 $\|D\|_1 = \|\tilde{D}\|_1 = 1/(1-\gamma)$ 和 $\|\Delta\|_1 = 2\alpha$, 依此得

$$\begin{aligned}
\|D\Delta\tilde{D}\chi\|_1 &\leq \\
\|D\|_1 \|\Delta\|_1 \|\tilde{D}\|_1 \|\chi\|_1 &= \\
= \frac{1}{1-\gamma} \cdot 2\alpha \cdot \frac{1}{1-\gamma} \cdot 1 &\quad (39)
\end{aligned}$$

因此, 最终证得:

$$\begin{aligned}
\gamma^2 |rD\Delta\tilde{D}\chi| &\leq \gamma \| \gamma r D \Delta \|_\infty \| D \Delta \tilde{D} \chi \|_1 \\
&\leq \gamma \| v \Delta \|_\infty \| D \Delta \tilde{D} \chi \|_1 \\
&\leq \gamma \| v \Delta \|_\infty \| D \Delta \tilde{D} \chi \|_1 \quad (40) \\
&\leq \gamma \cdot 2\alpha\rho \cdot \frac{2\alpha}{(1-\gamma)^2} = \frac{4\gamma\rho}{(1-\gamma)^2} \alpha^2
\end{aligned}$$

参考文献:

- [1] Wang D, Yin J, Tang T, et al. A two-step weighted least-squares method for joint estimation of source and sensor locations: A general framework [J]. Chinese Journal of Aeronautics, 2019, 32(2): 417-443.
- [2] 王鼎, 尹洁昕, 朱中梁. 针对超视距短波辐射源的测角与测时差协同定位方法[J]. 中国科学:信息科学, 2022, 52(11): 1942-1973.
Wang D, Yin J X, Zhu Z L. Novel cooperative localization method of over-the-horizon shortwave emitters based on direction-of-arrival and time-difference of-arrival measurements. Scientia Sinica Informationis, 2022, 52: 1942 - 1973. (in Chinese)
- [3] Zhou C, Gu Y, Shi Z, et al. Structured nyquist correlation reconstruction for DOA estimation with sparse arrays[J]. IEEE Transactions on Signal Processing, 2023, 71: 1849-1862.
- [4] Xia N, Xing B. A direct localization method for HF source geolocation and experimental results[J]. IEEE Antennas and Wireless Propagation Letters, 2021, 20(5): 728-732.
- [5] Zheng S, Yang Z, Shen W, et al. Deep learning-based DOA estimation [J]. IEEE Transactions on Cognitive Communications and Networking, 2024, 10(3): 819-835.
- [6] Adarsh M. S, Argyrios G, James Z., et al. Track-MDP: Reinforcement learning for target tracking with controlled sensing[C]// Proceedings of the 50th IEEE International Conference on Acoustics, Speech, and Signal Processing. Hyderabad, India: IEEE, 2025: 1-5.
- [7] Caicedo J C, Lazebnik S. Active object localization with deep reinforcement learning[C]// Proceedings of the 14th International Conference on Computer Vision. Santiago, Chile 2015: 2488-2496.
- [8] Bellver M, Giro X, Marques F, et al. Hierarchical object detection with deep reinforcement learning[C]// Proceedings of the 30th International Conference on Neural Information Processing Systems. Barcelona, Spain., 2016: 1-9.
- [9] Hao X, Yang S, Liu R, et al. VSLM: Virtual signal large model for few-shot wideband signal detection and recognition[J]. IEEE Transactions on Wireless Communications, 2025, 24(2): 909-925.
- [10] Li Y, Hu X, Zhuang Y, et al. Deep reinforcement learning (DRL): Another perspective for unsupervised wireless localization[J]. IEEE Internet of Things Journal, 2020, 7(7): 6279-6287.
- [11] Paul A, Singh K, Kaushik A, et al. Quantum-enhanced DRL optimization for DoA estimation and task offloading in ISAC systems[J]. IEEE Journal on Selected Areas in Communications, 2025, 43(1): 364-381.
- [12] Zhao L, Liu X, Shang T. Maximizing coverage in UAV-based emergency communication networks using deep reinforcement learning[J]. Signal Processing, 2025, 230: 109844.
- [13] Fu H, Tang H, Hao J, et al. Deep multi-agent reinforcement learning with discrete-continuous hybrid action spaces[C]// Proceedings of the 28th International Joint Conference on Artificial Intelligence. Macao, China, 2019: 2329-2335.
- [14] Upadhyay P, Marriboina V, Goyal S J, et al. An improved deep reinforcement learning routing technique for collision-free VANET[J]. Scientific Reports, 2023, 13(1): 1-12.
- [15] Xue W, Wu H, Ye H, et al. An improved proximal policy optimization method for low-level control of a quadrotor[J]. Actuators, 2022, 11(4): 105.
- [16] Wang R, Xu C, Sun J, et al. Cooperative localization for multi-agents based on reinforcement learning compensated filter[J]. IEEE Journal on Selected Areas in Communications, 2024, 42(10): 2820-2831.
- [17] Ran Wang, Jing Sun, Cheng Xu, et al. Reinforcement learning compensated filter for multi-agents cooperative localization[C]// Proceedings of the 49th IEEE International Conference on Acoustics, Speech, and Signal Processing.. Seoul, Korea: IEEE, 2024: 101-105.
- [18] Alagha A, Singh S, Mizouni R, et al. Target localization using multi-agent deep reinforcement learning with proximal policy optimization [J]. Future Generation Computer Systems, 2022, 136: 342-357.
- [19] 唐涛, 王鼎, 杨宾, 等. 一种基于高分辨测角语音图智能识别的短波自动测角方法: CN110045322B[P]. 2021.
Tang T, Wang D, Yang B, et al. A short-wave automatic direction finding method based on intelligent recognition of high-resolution direction finding time-frequency spectrum: CN110045322B[P]. 2019. (in Chinese)
- [20] 王鼎. 无线电测角与定位[M]. 国防工业出版社, 2016: 256-257.
Wang D, et al. Radio direction finding and localization [M]. National Defense Industry Press, 2016: 256-257. (in Chinese)
- [21] Schulman J, Levine S, Moritz P, et al. Trust region policy optimization [C]// Proceedings of the 32nd International Conference on Machine Learning. Lille, France. 2015: 1-16.
- [22] Gu Y, Cheng Y, Chen C L P, et al. Proximal policy optimization with policy feedback[J]. IEEE Transactions on Systems Man Cybernetics-Systems, 2022, 52(7): 4600-4610.

- [23] Lowe R, Wu Y, Tamar A, et al. Multi-agent actor-critic for mixed cooperative-competitive environments[J]. *Neural Information Processing Systems*, 2017: 1-16.
- [24] Iqbal S, Sha F. Actor-attention-critic for multi-agent reinforcement learning[C]//*Proceedings of the 36th International Conference on Machine Learning*. California, USA. 2019: 1-14.
- [25] Hu X, Li J, Zhan X, et al. Query-policy misalignment in preference-based reinforcement learning[C]// *Proceedings of the 12th International Conference on Learning Representations*. Vienna Austria. 2024: 1-25.
- [26] Zhu Y, Liu J, Gu P, et al. Improving reward models with proximal policy exploration for preference-based reinforcement learning[C]// *The 39th Annual Conference on Neural Information Processing Systems*. San Diego, USA. 2025: 1-13.
- [27] Yang Y, Zheng J, Liu H, et al. Optimal sensor placement and velocity configuration for TDOA - FDOA localization and tracking of a moving source[J]. *IEEE Transactions on Aerospace and Electronic Systems*, 2024, 60(6): 8255-8272.
- [28] Yang D, Solihin M I, Zhao Y, et al. Model compression for real-time object detection using rigorous gradation pruning[J]. *iScience*, 2025, 28(1): 111618.
- [29] 张旭辉, 姜春华, 刘桐辛, 等. 电离层虚高对超视距雷达多站联合定位精度的影响[J]. *电波科学学报*, 2022, 37(5) :761-767.
- Zhang X H, Jiang C H, Liu T X, et al. Effect of the ionospheric virtual height on the joint positioning accuracy of multi-station over-the-horizon radar system[J]. *Chinese journal of radio science*, 2022, 37(5): 761-767. (in Chinese).
- [30] Xia N, Xing B. A direct localization method for HF source geolocation and experimental results[J]. *IEEE Antennas and Wireless Propagation Letters*, 2021, 20(5): 728-732.



冯祺玥 (2002-), 女, 山东菏泽人, 信息工程大学博士生, 主要研究方向为强化学习、阵列信号处理和无源定位。



唐涛 (1981-), 男, 湖北荆门人, 博士, 信息工程大学教授, 博士生导师, 主要研究方向为阵列信号处理和无源定位。