

---

# 智能空中博弈算法挑战赛

## 智能体开发指南

南开大学  
二〇二四年七月

---

# 目 录

1 智能体开发流程指导说明 .....	3
1.1 概述 .....	3
1.2 博弈推演平台构建与配置 .....	3
2 智能体训练流程 .....	5
2.1 环境设计 .....	5
2.1.1 状态编码 .....	5
2.1.2 动作处理 .....	7
2.1.3 回报设计 .....	8
2.2 模型架构 .....	10
2.3 算法更新 .....	12
2.3.1 值网络损失 .....	13
2.3.2 策略网络损失 .....	13
2.4 训练主体 .....	13
3 智能体提交说明 .....	15
3.1 基类智能体 .....	15
3.2 自定义智能体 .....	16

---

## 1 智能体开发流程指导说明

### 1.1 概述

人工智能开发接口提供使用 Python 脚本编写的推演交互功能模块，Python 程序主要包含两个部分：环境模块（custom\_env 文件夹）、训练主体类模块（trainable\_class）与智能体模块（agent 文件夹）。param.py 包含了智能体训练所需用到的参数，main\_train.py 包含了训练的函数调用示例。

如图 1 所示，agent 文件夹下包含基类智能体（base\_agent.py）和红蓝双方规则智能体（red\_agent.py 与 blue\_agent.py）的示例程序；custom\_env 文件夹下包含基类环境类（env\_zhikong.py）和 1v1 示例环境类（DogFight.py）；训练主体类模块包含自定义模型文件夹（custom\_model）与自定义策略文件夹（custom\_policy）。

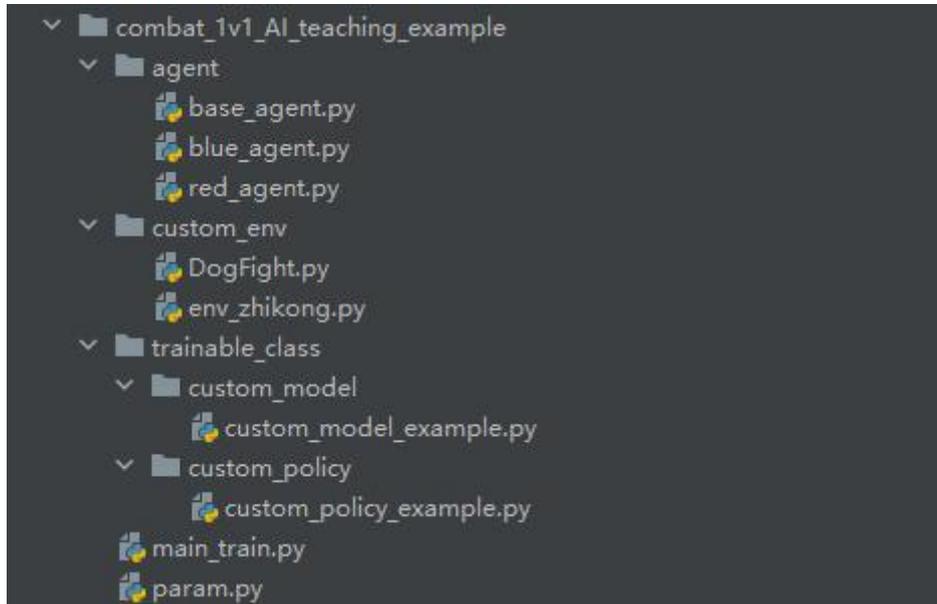


图 1 Python 程序架构

### 1.2 博弈推演平台构建与配置

仿真后台（仿真引擎）可以通过命令行直接启动，算法定义在 custom\_env/env\_zhikong.py 文件中，具体启动代码如图 2 所示。

```
def create_entity(self):
    is_success = False
    while not is_success:
        try:
            self.excute_cmd = f'{self.excute_path} Ip={self.IP} Port={self.PORT} ' \
                              f'PlayMode={self.RENDER} ' \
                              f'RedNum={self.red_num} BlueNum={self.blue_num} ' \
                              f'Scenes={self.scenes}'
            print('Creating Env', self.excute_cmd)
            self.unity = os.popen(self.excute_cmd)
            time.sleep(20)
            self._connect()
            is_success = True
            print('Env Created')
        except Exception as e:
            print('Create failed and the reason is ', e)
            time.sleep(5)
```

图2 仿真引擎启动

启动参数可在环境中自行定义，具体参数与含义简介如下：

Ip=ip 地址（127.0.0.1）

Port=端口号（8888）

RedNum=红方数量

BlueNum=蓝方数量

PlayMode=0/1（训练模式/演示模式）

Scenes=1/2/3/4/5（阿尔卑斯山/阿塔卡马沙漠/空场景/黄石公园/科罗拉多大峡谷）

若不输入可为默认值 Ip 127.0.0.1 Port 8888 RedNum 6 BlueNum 6  
PlayMode 0 Scenes=1。

---

## 2 智能体训练流程

为了便于开发使用，本部分以 1v1 为例，简要阐述开发流程。

RLlib 是开源的工业级强化学习库，可高效分配训练资源（CPU、GPU 与内存等），并提供优秀的异常处理等功能，可为算法的训练效率提供保障；近端裁剪策略梯度算法（PPO）是当前常用的强化学习基线算法，并以训练稳定著称，因此本实验基于 RLlib 实现工程开发，以 PPO 为底层更新算法，训练得到 1v1 场景下的固定翼无人机博弈智能体。

算法主体由环境设计、模型架构、算法更新与训练主体四部分组成。具体而言，环境设计部分实现了状态编码、动作处理与回报设计，对博弈问题进行建模并编码；模型架构部分设计了算法训练所需的深度神经网络，本实验包括值函数网络与策略函数网络；算法更新部分定义了模型优化器与损失函数，用于智能体的更新；训练主体则负责环境设计、模型架构与算法更新三部分的调配，实现多线程间的协调与资源分配，为智能体采样与更新的有序进行提供保障。

### 2.1 环境设计

环境设计部分由状态编码、动作处理与回报设计三部分组成，本小节详细阐述了其具体含义并提供具体代码样例，以方便开发者的后续开发，具体可见 `custom_env/DogFight.py`。

#### 2.1.1 状态编码

状态  $\mathbf{s}$  是智能体对自身与环境状态的数学编码，本实验假设环境为完全可观，为此可获取敌方任意时刻的位姿与武器信息，为我方智能体决策提供辅助。

---

本实验设计的状态编码包括红蓝双方

- (1) 基础信息：生命值与高度共 2 维；
- (2) 位姿信息：经纬度与姿态共 6 维；
- (3) 速度信息：机体坐标系下三个方向的速度信息共 3 维；
- (4) 角速度信息：机体坐标系下三个方向的角速度信息共 3 维；
- (5) 武器信息：两种武器的剩余数量、是否加载完毕与是否被地方武器锁定等信息共 9 维。

红蓝双方各包含 23 维单体状态，加入 2 维死亡信息，共计 48 维。  
为阐述更加清晰，单体状态编码部分的代码表示如下：

```
1. 单体无人机的状态编码
2. def _state_process(self, flag):
3.     """
4.     TODO: 本函数是状态编码信息，此处只是编码了生命值和高度信息
5.     TODO: 编程开发者可以参考智空文档，自行定义智能体训练需要的观测信息
6.     """
7.     if self.is_done[flag]:
8.         return [0 for i in range(23)]
9.     else:
10.        post_process_obs = []
11.        obs = self.obs_tot[flag][f'{flag}_0']
12.        # 生命值和高度
13.        post_process_obs.append(obs['LifeCurrent'])
14.        post_process_obs.append(obs['position/h-sl-ft'])
15.        # 经纬度
16.        post_process_obs.append(obs['position/long-gc-deg'])
17.        post_process_obs.append(obs['position/lat-geod-deg'])
18.        # 姿态
19.        post_process_obs.append(obs['attitude/pitch-rad'])
20.        post_process_obs.append(obs['attitude/roll-rad'])
21.        post_process_obs.append(obs['attitude/psi-deg'])
```

```

22.     post_process_obs.append(obs['acro/beta-deg'])
23.     # 速度
24.     post_process_obs.append(obs['velocities/u-fps'])
25.     post_process_obs.append(obs['velocities/v-fps'])
26.     post_process_obs.append(obs['velocities/w-fps'])
27.     # 角速度
28.     post_process_obs.append(obs['velocities/p-rad_sec'])
29.     post_process_obs.append(obs['velocities/q-rad_sec'])
30.     post_process_obs.append(obs['velocities/r-rad_sec'])
31.     # 武器
32.     post_process_obs.append(obs['SRAAMCurrentNum'])
33.     post_process_obs.append(obs['AMRAAMCurrentNum'])
34.     post_process_obs.append(obs['SRAAM1_CanReload'])
35.     post_process_obs.append(obs['SRAAM2_CanReload'])
36.     post_process_obs.append(obs['AMRAAM1_CanReload'])
37.     post_process_obs.append(obs['AMRAAM2_CanReload'])
38.     post_process_obs.append(obs['AMRAAM3_CanReload'])
39.     post_process_obs.append(obs['AMRAAM4_CanReload'])
40.     post_process_obs.append(obs['MissileAlert'])
41.     return post_process_obs

```

### 2.1.2 动作处理

动作 $\mathbf{a}$ 是智能体与环境交互时执行的指令，本实验中，智能体需在平稳控制无人机的同时，合理分配武器切换与攻击指令，因此，本实验的动作可表示为：

$$\mathbf{a} = [\mathbf{a}_m, \mathbf{a}_w],$$

其中 $\mathbf{a}_m, \mathbf{a}_w$ 分别表示机动控制指令与武器控制指令。

本实验通过油门、副翼、升降舵与方向舵实现对无人机的控制，即

$$\mathbf{a}_m = [c_e, c_a, c_r, c_T],$$

其中 $c_e, c_a, c_r, c_T$ 分别表示无人机的升降舵、副翼、方向舵与油门。

武器控制指令 $a_w$ 包括武器切换与武器发射，可表示如下：

$$a_w = [l_w, s_w].$$

值得注意的是，机动控制指令为连续控制量，武器控制指令则为离散控制量，为保持动作类型的一致性，对机动控制量进行离散化处理，结合武器控制指令，动作空间最终可表示如下：

$$\begin{cases} c_e \in \{-1, 0, 1\} \\ c_a \in \{-1, 0, 1\} \\ c_r \in \{-1, 0, 1\} \\ c_T \in \{0, 0.5, 1\} \\ l_w \in \{0, 1\} \\ s_w \in \{0, 1\} \end{cases}$$

综上，连续-离散混合动作空间维度为 324 维，动作解码部分代码表示如下：

```
1. 动作解码
2. def input_index_action(self, action):
3.     action = int(action)
4.     switch_missile_action = int(action // 162)
5.     action %= 162
6.     launch_missile_action = int(action // 81)
7.     action %= 81
8.     action_ce = action // 27
9.     action %= 27
10.    action_ca = action // 9
11.    action %= 9
12.    action_cr = action
13.    action %= 3
14.    action_cT = action
```

---

### 2.1.3 回报设计

回报 $r$ 是智能体与环境进行交互后，环境对智能体行为好坏的反馈，其设计方式直接影响着智能体的训练方向与训练速度。

在 1v1 博弈问题中，智能体旨在击杀对手的同时保证自身存活，一时的优势或劣势地位无法改变博弈的根本目标，基于威胁指数的稠密回报设计可能使智能体过度在意一时的得失，忽略博弈的根本目的，导致训练偏离博弈目标；然而仅使用胜负作为回报则会引起稀疏回报问题，这是由于 1v1 博弈问题具有交互序列长等特点，常常需要成百上千步才能得到最终的胜负信息。因此，如何权衡稠密回报与稀疏回报的设计，是留给开发者解决的一大难题。

为便于阐述，本实验使用基于胜负信息的稀疏回报，具体瞬时回报设计如下：

$$r = \begin{cases} 1, & \text{if win} \\ 0, & \text{if tie} \\ -1, & \text{if loss} \end{cases} .$$

回报设计的核心代码表示如下：

```
1. 瞬时回报设计
2. def get_win_loss_reward(self):
3.     """
4.     Returns: win_loss reward
5.     """
6.     control_side = self.control_side
7.     opponent_side = self.opponent_side
8.     if self.is_done['__all__']:
9.         control_side_death_event = self.obs_tot[control_side][f'{control_side}_0']['DeathEvent']
10.        opponent_side_death_event = self.obs_tot[opponent_side][f'{opponent_side}_0']['DeathEvent']
11.        # 1. 没打死对方给予负向回报
12.        if control_side_death_event == 99 and opponent_side_death_event == 99:
```

```

13.     return 0
14.     # 2. 红蓝都死了，激进型目的达到部分，武器给予正向回报
15.     if control_side_death_event != 99 and opponent_side_death_event != 99:
16.         return 0
17.     # 3. 红方死了，蓝方没死，一定给予惩罚
18.     if control_side_death_event != 99:
19.         return -1
20.     # 4. 红方没死，蓝方死了，给予奖励
21.     else:
22.         return 1
23.     return 0

```

## 2.2 模型架构

模型架构部分设计了算法训练所需的深度神经网络，具体可见 trainable\_class/custom\_model/custom\_model\_example.py 文件，本实验选用 PPO 算法，共包括值函数网络与策略函数网络，网络的示意图如图所示。

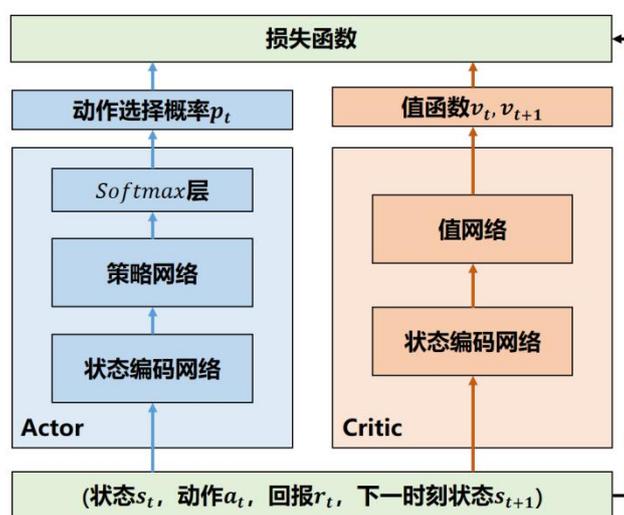


图3 网络设计示意图

为便于开发者学习，我们将网络定义的核心部分代码如下，开发者可根据项目需求自行更改网络设置：

```

1.  模型架构核心代码
2.  class Custom_Model_Example(TorchModelV2, nn.Module):
3.      def __init__(
4.          self, obs_space: gym.spaces.Space,
5.          action_space: gym.spaces.Space, num_outputs: int,
6.          model_config: ModelConfigDict, name: str,
7.      ):
8.          TorchModelV2.__init__(
9.              self, obs_space, action_space, num_outputs, model_config, name
10.         )
11.         nn.Module.__init__(self)
12.         hiddens = list(model_config.get("fcnet_hiddens", [])) + list(
13.             model_config.get("post_fcnet_hiddens", [])
14.         )
15.         activation = model_config.get("fcnet_activation")
16.         if not model_config.get("fcnet_hiddens", []):
17.             activation = model_config.get("post_fcnet_activation")
18.         #####
19.         # 策略网络的隐藏层
20.         #####
21.         layers = []
22.         prev_layer_size = int(np.product(obs_space.shape))
23.         self._logits = None
24.         for size in hiddens[:-1]:
25.             layers.append(
26.                 SlimFC(
27.                     in_size=prev_layer_size, out_size=size,
28.                     initializer=normc_initializer(1.0), activation_fn=activation
29.                 ))
30.             prev_layer_size = size
31.         if len(hiddens) > 0:
32.             layers.append(
33.                 SlimFC(
34.                     in_size=prev_layer_size, out_size=hiddens[-1],
35.                     initializer=normc_initializer(1.0), activation_fn=activation,

```

```

36.     ))
37.     prev_layer_size = hiddens[-1]
38.     self._logits = SlimFC(
39.         in_size=prev_layer_size, out_size=num_outputs,
40.         initializer=normc_initializer(0.01), activation_fn=None,
41.     )
42.     self._hidden_layers = nn.Sequential(*layers)
43.     #####
44.     # 值函数的网络部分，包括了值函数的编码网络与最后的输出层
45.     #####
46.     prev_vf_layer_size = int(np.product(obs_space.shape))
47.     vf_layers = []
48.     for size in hiddens:
49.         vf_layers.append(
50.             SlimFC(
51.                 in_size=prev_vf_layer_size, out_size=size,
52.                 activation_fn=activation, initializer=normc_initializer(1.0),
53.             ))
54.         prev_vf_layer_size = size
55.     self._value_branch_separate = nn.Sequential(*vf_layers)
56.     self._value_branch = SlimFC(
57.         in_size=prev_layer_size, out_size=1,
58.         initializer=normc_initializer(0.01), activation_fn=None,
59.     )
60.     self._features = None
61.     self._last_flat_in = None

```

## 2.3 算法更新

算法更新部分定义了模型优化器与损失函数，用于智能体的更新，具体可见 `trainable_class/custom_policy/custom_policy_example.py` 文件。模型优化器采用 Adam 优化器，同时优化值网络与策略网络。损失函数由值网络损失与策略网络损失组成，值网络损失与策略网络

---

损失的具体计算方法已在 2.2 中进行了详细阐述。

### 2.3.1 值网络损失

为便于开发者学习，我们将值网络损失核心代码展示如下：

```
1. 值网络损失计算
2. value_fn_out = model.value_function()
3. vf_loss = torch.pow(
4.     value_fn_out - train_batch[Postprocessing.VALUE_TARGETS], 2.0
5. )
6. vf_loss_clipped = torch.clamp(vf_loss, 0, self.config["vf_clip_param"])
7. mean_vf_loss = reduce_mean_valid(vf_loss_clipped)
```

### 2.3.2 策略网络损失

策略网络损失核心代码展示如下：

```
1. 策略网络损失计算
2. logp_ratio = torch.exp(
3.     curr_action_dist.logp(train_batch[SampleBatch.ACTIONS])
4.     - train_batch[SampleBatch.ACTION_LOGP]
5. )
6. surrogate_loss = torch.min(
7.     train_batch[Postprocessing.ADVANTAGES] * logp_ratio,
8.     train_batch[Postprocessing.ADVANTAGES]
9.     * torch.clamp(
10.         logp_ratio, 1 - self.config["clip_param"], 1 + self.config["clip_param"]
11.     ),
12. )
13. mean_policy_loss = reduce_mean_valid(-surrogate_loss)
```

## 2.4 训练主体

训练主体类负责环境设计、模型架构与算法更新三部分的调配，实现多线程间的协调与资源分配，为智能体采样与更新的有序进行提

供保障。本实验选用的 PPO 算法是典型 on-policy 算法,因此遵循“采样-更新”迭代更新架构。具体而言:在采样阶段,训练主体类接收环境返回的状态,根据预设探索策略与模型计算当前时刻动作,与环境进行交互并获得反馈信息;待采样足够数据后,进入更新阶段,此时训练主体类从缓存中采集训练样本,根据损失更新采样模型。整个过程中,训练主体类负责保证采样与更新阶段的模型参数的实时性,为智能体采样与更新的有序进行提供保障。

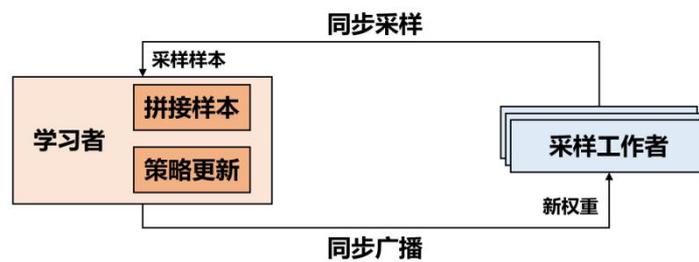


图 4 训练流程

为便于开发者学习,我们将训练主体类的组成框架部分代码展示如下,开发者可在 `training_step(self)` 函数中定义采样与训练的时序关系(同步/异步等),具体方法可参考 RLlib 教学文档<sup>[1]</sup>:

```

1. class Custom_Trainable_Policy(PPO):
2.     @override(PPO)
3.     def get_default_policy_class(self, config):
4.
5.     @classmethod
6.     @override(Algorithm)
7.     def get_default_config(cls) -> AlgorithmConfigDict:
8.
9.     @override(Algorithm)
10.    def validate_config(self, config: AlgorithmConfigDict) -> None:
11.
12.    @ExperimentalAPI

```

---

13. `def training_step(self) -> ResultDict:`

### 3 智能体提交说明

#### 3.1 基类智能体

为了方便统一对战调用形式，在文件夹 `agent/base_agent.py` 中定义了一个基类智能体：`Base_agent`。该基类包含 `get_action` 函数，用于接受平台态势数据，返回智能控制指令。平台返回的态势数据与返回平台的控制指令的数据格式详见文档“\*\*\*\*”。

`Base_Agent` 类的具体属性与具体方法如表 1 与表 2 所示。

Base\_Agent 类的属性

属性	类型	说明
<code>side</code>	<code>string</code>	智能体对应的推演方
<code>control_num</code>	<code>int</code>	智能体包含的智能体数量
<code>agent_config</code>	<code>dict</code>	智能体运行所需的额外参数

Base\_Agent 类的方法功能

方法	功能说明
<code>__init__</code>	智能体初始化
<code>agent_step</code>	智能体决策
<code>reset</code>	智能体重置

`Base_Agent` 的具体代码定义如图 5 所示。

```

1 class Base_agent:
2     def __init__(self, control_num, side, agent_config):
3         """
4         :param control_num: 我方需要控制的智能体数量
5         :param side: 我方控制的阵营, 包含'red'与'blue'
6         """
7         self.side = side
8         self.control_num = control_num
9         self.agent_config = agent_config
10        self.obs_tot = None
11
12    def reset(self):
13        pass
14
15    def agent_step(self, obs):
16        """
17        :param obs: 平台提供的返回信息
18        :return: 需要向平台发送的操作信息
19        """
20        pass

```

图 5 Base\_Agent 基类智能体定义

### 3.2 自定义智能体

选手自定义的智能体必须继承基类智能体, 并在初始化时调用父类的初始化函数。

可按下述顺序完成智能体类: 态势信息构建和预处理—>搭建神经网络结构—>调用并加载已训练模型—>根据编码的态势生成动作—>按照平台接收数据格式规范化动作。

需要注意的是, 某一方智能体只能获取本推演方的态势信息(战损信息除外), 也只能对本推演方的平台进行控制, 目前的发布版本并没有进行相关限制, 但比赛进行时仿真引擎会检查 python 的命令有效性并过滤非法指令(对仿真逻辑没有影响)。