

Contents

1	Dr. Yielding's Questions	1
2	Dr. Robbins' Questions	2
	Multi-Agent Lunar Lander	3
3	Dr. Cox's Questions	6
A	Experiment Running Script	7
B	Multi-Agent Lunar Lander	12

1 Dr. Yielding's Questions

While HARL is indeed a subfield of MARL, the distinction between the two can often be unclear or inconsistently used in the literature. In the broadest sense, 'heterogeneous' may refer to any MARL scenario where the agents are not identical. However, using this broad definition risks diluting the practical significance of the term.

To provide a more useful distinction, HARL should be invoked when the heterogeneity of the agents is fundamental, essential, or definitional rather than merely incidental. This means the differences among agents are critical to their roles and interactions within the system, not just minor variations.

Prior to my literature review, I would have considered HARL to apply specifically to cases where agents are distinct from the outset, either in their capabilities (For action-space \mathcal{A} , $\mathcal{A}_1 \neq \mathcal{A}_2$) or their observation spaces ($\mathcal{O}_1 \neq \mathcal{O}_2$). This intrinsic heterogeneity is evident in works like [1] (the Irish conference paper mentioned in Dr. Yielding's question) and implicitly reflected in [2] (OpenAI Five), even though they do not explicitly label their methods as HARL.

The most comprehensive source on HARL usage is Zhong et al. [3], whose work has greatly influenced my understanding. They focus on implementing algorithms that encourage the development of heterogeneous policies among agents. Their framework increases the likelihood of individual agents converging on distinct policies, which I term emergent heterogeneity.

In my prospectus, I also mention a suspicion that this approach might not entirely prevent agents from converging on policies that are functionally similar, thus lacking true diversity. However, this assertion remains an ancillary detail as it is not yet substantiated by empirical evidence.

Therefore, I propose distinguishing between intrinsic heterogeneity, where agents are fundamentally different before training, and emergent heterogeneity, where differences arise as a result of the learning process. In the cases where the heterogeneity of the agents falls below a level of functional relevance, and appears to be incidental, I argue that they should not be labeled as HARL. By maintaining these distinctions, we can more accurately categorize and understand the applications and implications of HARL and MARL.

Below, I apply these distinctions to the cases proposed:

1.1 HAA2C and HADDPG are among the numerous algorithms proposed by Zhong et al. [3]. While it seems reasonable that their algorithms could be applied to situations with intrinsic heterogeneity, their implementations and tests are applied to environments with agents that are functionally the same.

I intend to implement (at least a subset of) their algorithms to the extent that time allows. In doing so, I hope to either corroborate or contradict their results by comparing them to similar algorithms under different

conditions. This exploration aims to identify any apparent advantages of these algorithms when applied to the experimental variables proposed for Contribution 1. These experimental variables represent smaller difficulties that we expect to face in Contributions 2 and 3.

- 1.2 (a) The scenario described in this part of the question is, perhaps unintuitively, more akin to single-agent than multi-agent reinforcement learning. This becomes clearer when considering a single agent acting as an 'overlord,' where the observation space is a combination of observations from all the individual agents. The actions taken by this overlord are combinations of actions chosen for each agent. Essentially, a single policy processes the combined observation and outputs the combined action.
- (b) This example is a strong example of MARL, and because the agents utilize copies of a singular policy, this example is free from any of the types of heterogeneity described in the answer for question 1.1.
- (c) The types of problems accurately described by the scenario provided in this part of the question appear to be a subset of MARL problems and a superset of HARL problems.

Many MARL algorithms allow the member policies to develop distinctly (e.g., [4]–[6]), but they are not optimized to facilitate the development of distinct policies.

Zhong et al. [3] formulate their series of HARL algorithms with optimizations intended to facilitate the development of distinct policies. One weakness of this formulation is that there is no guarantee that the multiple policies will not converge to a behaviorally indistinct set, similar to the concept of carcinization observed in evolutionary biology.

Thus, the resulting heterogeneity of the agents in these scenarios is not intrinsic but emergent. Whether the algorithm itself is labeled as MARL or HARL is distinguished by intent.

- 1.3 Referencing Centralized Training Decentralized Execution (CTDE) and Decentralized Training Decentralized Execution (DTDE) as employed by Li et al. in their FA2A paper [7], we see that CTDE is the most common format for Actor-Critic based MARL algorithms [4]–[8].

Li et al. [7] and Wen et al. [9] are the only papers I found that discuss the contrary implementation, DTDE. In both cases, the authors motivate DTDE with practical concerns, particularly the limitations of inter-agent communication in distributed systems. These considerations are important, but the relation to HARL remains the same as described in answer 1.2 (c).

2 Dr. Robbins' Questions

In addition to being included in the appendices of this exam, all code used to run the experiments is available on my github at <https://github.com/bhosley/Specialty-Exam>. It is written to be run on any ANT-Center VSCode server containers, but should work in a generic virtual environment. The specifications for the virtual machines used for this experiment are enumerated in table 1.

I recommended increasing the ratio of memory to CPU allocation for future experiments as the container was substantially closer to maximum utilization of memory than processing for the duration of the experiments run for this exam.

The github readme has a short guide for setting up the virtual environment in the ANT-Center for easy replication. The answers for this section are drawn from the experiment running script, also provided in appendix A of this document.

Setting	Value
Template	vscode-server
Image	reg.git.act3-ace.com/ace/hub/vscode-server:v0
Max CPUs	64
Req. CPUs	16
Memory	64 GB
GPU	None

Table 1: Container Settings

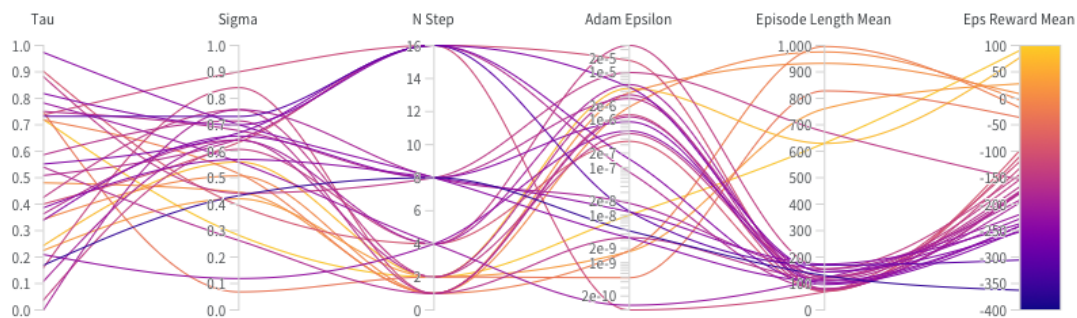


Figure 1: Lunar Lander Baseline Experiment 1 Hyperparameters

2.1 The Deep Q-network (DQN) implementation [10], [11] implemented in the RLlib framework [12] was used for each of the answers in this section. The default results output is readable by tensorboard, but I chose to use the wandb (Weights and Balances) api as well.

To perform this baseline experiment we can call the default script (appendix A). We can accomplish a 5 replication test by setting `--num-env-runners=5`, however, this overrides the default of 10, so we have chosen not to use it and to keep the default. We use the

```
python dqn_exp.py --sweep --num-samples=30 --num-env-runners=60
```

Round 1 - Bad Results:

In the first round of Parameter sweeping, a large number of training runs failed to converge on reasonable results. This was somewhat surprising, and when examining the parameters there was no immediately obvious pattern. Adding the average number of steps per episode to the Comparison (fig. 1) shows consistently shorter episodes for poor performance. Further investigation showed that the observation space assigned to the environment from the RLlib registered lunar lander environment was significantly different from the correct space. Further, manually registering the environment after importing it directly from the Farama maintained packages using Ray's environment registration function resulted in an outdated observation space, less prone to resulting in error, but still problematic.

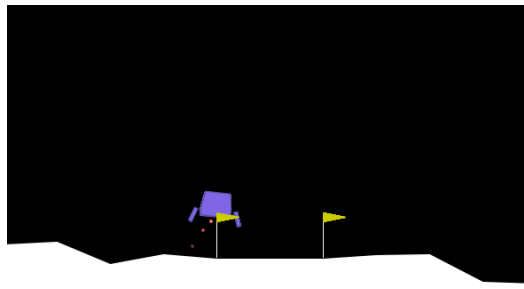
Thus we gain two actionable insights. First, we adjust this experiment to use the custom Lunar Lander environment described in the next section, but set to one agent. Second, we identify an update that can be performed by a pull request to the maintainers of RLlib, a contribution to the open-source software I would like to attempt after the submission of this exam.

Round 2 - Baseline Parameter Sweep:

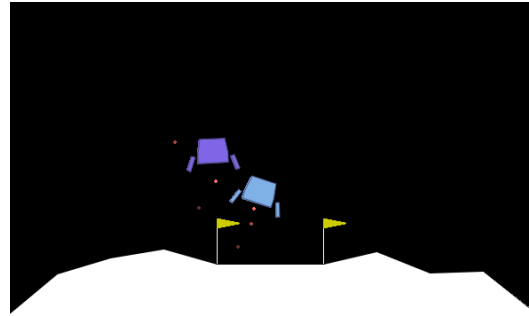
Multi-Agent Lunar Lander

To modify the Lunar Lander environment to support multiple agents it appeared to be possible to simply duplicate all of the sections that referenced the lander in the original code. However, to achieve cleaner code, and behavior more consistent with Farama's PettingZoo environments, I elected to rewrite the environment to be more consistent with object oriented programming principles.

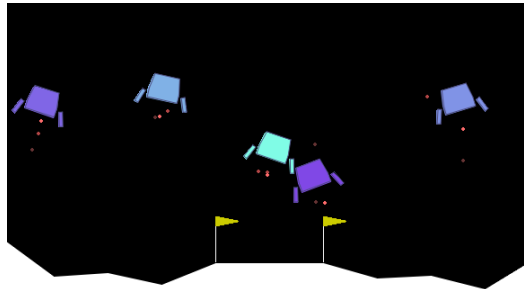
This trivialized the retention of the environmental physics, and made it very easy for the new environment to comply with both the AEC (Agent Environment Cycle; or sequential) and parallel execution methods used in the



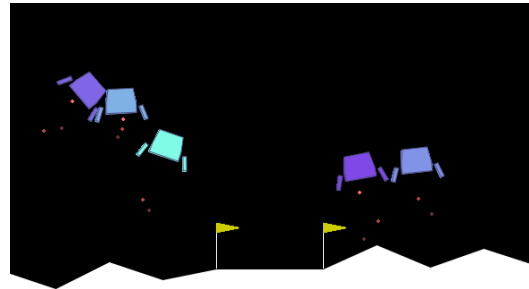
(a) Original



(b) With 2 Landers



(c) With 5 Landers



(d) 5 Landers with Collision

Figure 2: Multi-Agent Lunar Lander

PettingZoo API. Further, this made it much easier to retain the heuristic solution from the original environment. Figure 2 shows several screen shots from the resulting environment when rendered.

While an AEC version of this custom environment was constructed, it was not used for this exam. AEC iterates through the list of agents, allowing the agent to select an action and steps the environment with that action before moving to the next agent, in a manner similar to a board game. Such an implement doesn't necessarily contribute value to the multi agent Lunar Lander sim, and in fact, causes problems if the physics of the environment are not adjusted as the number of landers is increased. Using the heuristic as a measure, the AEC environment becomes unsolvable with 3 or more landers. The effect of their actions must be scaled to overcome the wait time between their turns, otherwise they simple crash, unable to overcome gravity.

The parallel version of the environment was used for this exam, and required only one major assumption, centering around collisions between the landers. Lunar Lander uses the Box2D package to model the moon surface and lander. The most expedient method to address collisions using Box2D was to detect only contact with the lander's body. If any other object touches a lander's body it is considered a crash. As written, the Box2D contact detection does not distinguish between what the secondary object is, if the legs of the lander where included, contact with lunar surface would show as a crash. Thus I chose not to include the legs contact detection for collisions. The effective result is that if two lander's legs touch it is not treated as a crash, however, if one lander's legs touch another's body module it is treated as a crash and thus a failure.

2.2 Formulate a Markov Game: The Markov game that represents the multi-agent version of this environment is an extension of the Markov Decision Process (MDP) that represents the original Lunar Lander environment. The formulation presented here will be used to describe the interactions in all of the questions that follow.

Agent: The agent is represented as a member of set of agents $n \in N$.

State Space: Let \mathcal{S} be the state space. Then, $\mathcal{S} \equiv s^N$, where

$$s = \begin{cases} x \in [-2.5, 2.5] & \text{Position of } n \text{ in } x \\ y \in [-2.5, 2.5] & \text{Position of } n \text{ in } y \\ \vec{x} \in [-10, 10] & \text{Velocity of } n \text{ in } x \\ \vec{y} \in [-10, 10] & \text{Velocity of } n \text{ in } y \\ \omega \in [-2\pi, 2\pi] & \text{Angle of } n \\ \vec{\omega} \in [-10, 10] & \text{Angular Velocity of } n \\ \mathbb{I}(\text{leg } 1) \in \{0, 1\} & \text{Leg on ground} \\ \mathbb{I}(\text{leg } 2) \in \{0, 1\} & \text{Leg on ground} \end{cases}$$

Action Space: Let \mathcal{A} be the action space. Then, $\mathcal{A} \equiv a^N$, where

$$a \in \begin{cases} 0 & \text{No-Op} \\ 1 & \text{Fire Left Engine} \\ 2 & \text{Fire Main Engine} \\ 3 & \text{Fire Right Engine} \end{cases}$$

Transition Probability: The transition probability for each agent n remains the same as the original environment,

$$P_n(s_{n,t+1}|s_{n,t}, a_{n,t}) \cong \begin{cases} \text{Dispersion} \sim U(-1, 1) \\ \text{Wind(Linear)} = \tanh(\sin(2kx) + \sin(\pi kx)) \\ \text{Wind(Rotate)} = \tanh(\sin(2kx) + \sin(\pi kx)) \end{cases}$$

Thus the transition probability for the game as a whole is

$$P(s_{t+1}|s_t, a_t) = \prod_{n \in N} P_n$$

Reward: The reward structure is similarly retained from the original Lunar Lander environment. Specifically, the reward at time t is defined as $r(t) = \sum_{n \in N} r_n(t)$ where

$$r_n(t) = \begin{array}{ll} \pm 100 & \text{End-state, crash or land} \\ + 10(s_{t,6} + s_{t,7}) & \text{leg(s) on ground} \\ - a_n(t) \cdot [0, 0.03, 0.3, 0.03] & \text{thruster cost} \\ - 100\sqrt{s_n(t)_0^2 + s_n(t)_1^2} & \text{Distance} \\ - 100\sqrt{s_n(t)_2^2 + s_n(t)_3^2} & \text{Velocity} \\ - 100|\omega_n(t)| & \text{Tilt} \end{array}$$

2.3 DQN Single-Agent Controller: This problem was approached using a wrapper class around the custom multi-agent environment to translate the between the multi-agent environment and single-agent policy. The proxy single-agent state/observation was formed by simply concatenating the state vectors of each agent.

The modified action space for this problem is $\mathcal{B} \equiv \mathcal{A}^N$. To relate the two action spaces for this problem I use the function,

$$b = \sum a_n \times \dim(a)^n.$$

Then to translate the modified action into the agent-action vector necessary for the multi-agent environment I use the following function,

$$\mathbf{a} = \{b / \dim(a)^n \mod \dim(a)\}_{n \in N}$$

which relies on $\dim(\mathcal{A}_n) = \dim(\mathcal{A}_m) \forall n, m \in N$, an assumption that I note as it is true for this problem, limits some generality. The policy can thus be represented as $\pi(b | [s_n]_{n \in N})$.

Finally, the experiment can be replicated using the command:

```
python dqn_exp.py --SA --sweep --num-samples=30 --num-env-runners=30
```

2.4 NoPS - No Parameter Sharing: This experiment uses a distinct policy for each agent, $\pi_n(a_n | s_n)$, and can be replicated using:

```
python src/dqn_exp.py --NoPS --sweep --num-samples=30 --num-env-runners=30
```

2.5 FuPS - Full Parameter Sharing: This experiment uses a shared policy for each agent, $\pi(a_n | s_n)$, and can be replicated using:

```
python src/dqn_exp.py --FuPS --sweep --num-samples=30 --num-env-runners=30
```

2.6 Results Comparisons:

3 Dr. Cox's Questions

A Experiment Running Script

```
1  """Script for Dr. Robbins Questions
2  Establishing a baseline implementation of Single Agent Lunar Lander
3  with a DQN.
4
5  See: https://gymnasium.farama.org/environments/box2d/lunar_lander/
6  for more details on the environment.
7
8  How to run this script
9  -----
10 `python [script_name].py`
11
12 For debugging, use the following additional command line options
13 `--no-tune --num-env-runners=0`
14
15 For logging to your WandB account, use:
16 `--wandb-key=[your WandB API key] --wandb-project=[some project name]
17 --wandb-run-name=[optional: WandB run name (within the defined project)]`
18
19
20
21
22
23
24 python dqn_exp.py \
25     --num-env-runners=5 \
26     --stop-reward=200 \
27     --wandb-key=913528a8e92bf601b6eb055a459bcc89130c7f5f \
28     --wandb-project=lunar_lander_test
29
30 tensorboard --logdir logs/fit
31 """
32 # --evaluation-duration=5 \ # Default uses 10
33 # --checkpoint-at-end
34 # --num-samples=30
35 # result_dict['env_runners']['episode_return_mean']
36
37
38 import numpy as np
39 from gymnasium import Wrapper, spaces
40 from argparse import ArgumentParser
41
42 from ray import tune
43 from ray.rllib.core.rl_module.rl_module import SingleAgentRLModuleSpec
44 from ray.rllib.core.rl_module.marl_module import MultiAgentRLModuleSpec
45 from ray.rllib.env.wrappers.pettingzoo_env import ParallelPettingZooEnv
46 from ray.rllib.utils.test_utils import (add_rllib_example_script_args,
47                                           run_rllib_example_script_experiment)
48 from ray.tune.registry import get_trainable_cls, register_env
49
50 import multi_lander
51
52
53 parser = add_rllib_example_script_args(
54     ArgumentParser(conflict_handler='resolve'), # Resolve for num_agents
55     default_reward=0.0,
56     default_iters=100, #100, #200
57     default_timesteps=100000, #10000, #100000
58 )
```

```

59 parser.add_argument(
60     "--control",
61     type=str,
62     choices=["Baseline", "SA", "NoPS", "FuPS"],
63     default="Baseline",
64     help="The controller method."
65     "`Baseline`: Original Lunar Lander with single agent and lander"
66     "`SA`: A single agent controls all of the landers."
67     "`NoPS`: No parameter sharing between the agents."
68     "`FuPS`: Full parameter sharing between the agents.",
69 )
70 parser.add_argument(
71     "--SA", action="store_true",
72     help="`SA`: Creates a single agent for controlling all of the agents.",
73 )
74 parser.add_argument(
75     "--NoPS", action="store_true",
76     help="`NoPS`: No parameter sharing between the agents.",
77 )
78 parser.add_argument(
79     "--FuPS", action="store_true",
80     help="`FuPS`: Full parameter sharing between the agents.",
81 )
82 parser.add_argument(
83     "--num-agents", type=int, default=2,
84     help="The number of agents",
85 )
86 parser.add_argument(
87     "--sweep", action="store_true",
88     help="Perform a parameter sweep instead of using tune",
89 )
90
91
92 class CustomWrapper(ParallelPettingZooEnv):
93     """Wraps a Parallel Petting Zoo Environment for RLlib
94
95     This wrapper is necessary to interface with rllib's workers.
96     There is a problem caused when the one agent terminates before
97     the other, which seems to pass inconsistent length episode
98     trajectories to the rollout worker(s).
99
100     It also adds the necessary `__all__` key to the done dictionaries.
101     """
102     def step(self, action_dict) -> tuple:
103         obs, rew, terminateds, truncateds, info = self.par_env.step(action_dict)
104
105         active = [a_id for a_id, term in terminateds.items() if term==False]
106
107         obs_s = {a_id: obs.get(a_id) for a_id in active}
108         rew_s = {a_id: rew.get(a_id) for a_id in active}
109         terminateds = {a_id: terminateds.get(a_id) for a_id in active}
110         truncateds = {a_id: truncateds.get(a_id) for a_id in active}
111
112         terminateds["__all__"] = all(terminateds.values())
113         truncateds["__all__"] = all(truncateds.values())
114         return obs_s, rew_s, terminateds, truncateds, info
115
116 # Registry is necessary for functional passing later.
117 register_env("ma-lander", lambda _: CustomWrapper(multi_lander.Parallel_Env()))
118

```



```

119
120 class SingleAgentWrapper(Wrapper):
121     """Wraps a parallel multi-agent environment for single-agent control.
122
123     *It is currently limited to agents with discrete action spaces and
124 box observation spaces with identical lengths.
125
126     This will wrap environments with an arbitrary number of agents,
127 however, the action space scales exponentially by the number of
128 agents and is likely to be the limiting factor for performance.
129     """
130     def __init__(self, *args):
131         super().__init__(*args)
132         # Concat all action spaces into a single action space
133         self.action_space = spaces.Discrete(np.prod(
134             [list(self.env.action_spaces.values())[i].n
135              for i in range(len(self.env.observation_spaces))]))
136         self._act_n = list(self.env.action_spaces.values())[0].n
137         # Concat all obs spaces into a single obs space
138         low, high = [], []
139         for i in range(len(self.env.observation_spaces)):
140             low.extend(list(self.env.observation_spaces.values())[i].low)
141             high.extend(list(self.env.observation_spaces.values())[i].high)
142         self.observation_space = spaces.Box(np.array(low), np.array(high))
143
144     def observe(self) -> tuple:
145         # Concat all obs values into a single obs tuple
146         #return sum(self.env.observe().values(),[]) # By monoid
147         return np.array(list(self.env.observe().values())).flatten()
148
149     def observation(self) -> tuple:
150         return self.observe()
151
152     def reset(self, *, seed:int | None=None, options:dict | None=None) -> tuple:
153         self.env.reset(seed=seed, options=options)
154         return self.observe(), {}
155
156     def step(self, action:int) -> tuple:
157         action_list = {a: divmod(int(action/(self._act_n*i)),self._act_n)[1]
158                        for i,a in enumerate(self.agents)}
159         obss, rews, terminations, _, _ = self.env.step(action_list)
160         obs = self.observe()
161         reward = sum(rews.values())
162         term = all(terminations.values()) or self.env.game_over
163         return obs, reward, term, False, {}
164
165     # Registry is necessary for functional passing later.
166     register_env("sa-lander",
167         lambda _: SingleAgentWrapper(multi_lander.Parallel_Env()))
168
169
170 if __name__ == "__main__":
171     args = parser.parse_args()
172
173     # Shared config settings for all experiments
174     base_config = (
175         get_trainable_cls("DQN")
176         .get_default_config()
177         .framework('torch')
178     )

```

```

179
180 # Set config for experiment if using Single Agent Control
181 if args.control == 'NoPS' or args.SA:
182     config = (
183         base_config
184         .environment(
185             "sa-lander",
186             env_config={"num_landers": args.num_agents}
187         )
188     )
189 # Set config for experiment if using No-Parameter Sharing
190 elif args.control == 'NoPS' or args.NoPS:
191     policies = {"lander_" + str(i) for i in range(args.num_agents)}
192     config = (
193         base_config
194         .multi_agent(
195             policies=policies, # 1:1 map from AgentID to ModuleID.
196             policy_mapping_fn=(lambda aid, *args, **kwargs: aid),
197         )
198         .rl_module(
199             rl_module_spec=MultiAgentRLModuleSpec(
200                 module_specs={p:SingleAgentRLModuleSpec() for p in policies},
201             ),
202         )
203         .environment(
204             "ma-lander",
205             env_config={"num_landers": args.num_agents}
206         )
207     )
208 # Set config for experiment if using Full-Parameter Sharing
209 elif args.control == 'FuPS' or args.FuPS:
210     config = (
211         base_config
212         .multi_agent(
213             policies={"p0"}, # All agents map to the same policy.
214             policy_mapping_fn=(lambda aid, *args, **kwargs: "p0"),
215         )
216         .rl_module(
217             rl_module_spec=MultiAgentRLModuleSpec(
218                 module_specs={"p0": SingleAgentRLModuleSpec()},
219             ),
220         )
221         .environment(
222             "ma-lander",
223             env_config={"num_landers": args.num_agents}
224         )
225     )
226 # Default, original lunar lander with one agent and lander
227 else:
228     config = (
229         base_config
230         .environment(
231             #env="LunarLander-v2"
232             "ma-lander",
233             env_config={"num_landers": 1}
234         )
235     )
236
237 # Parameter sweep settings
238 if args.sweep:

```

```

239 param_space = {
240     "adam_epsilon": tune.loguniform(1e-4, 1e-10), # 1e-8
241     "sigma0": tune.randn(0.5, 0.2), # 0.5
242     "n_step": tune.choice([2**i for i in range(5)]), # 1
243     # Update the target by \tau * policy + (1-\tau) * target_policy.
244     "tau": tune.uniform(0.0,1.0), # 1.0,
245     #epsilon = [(0, 1.0), (10000, 0.05)] [(step,epsilon),...]
246     # -> 1.0 at beginning, decreases to 0.05 over 10k steps
247 }
248 config = config.training(**param_space)
249 # Use Param sweep, not tune
250 #args.no_tune = True
251
252 # Call experiment runner
253 run_rllib_example_script_experiment(config, args)

```

B Multi-Agent Lunar Lander

```
1 import math
2 from typing import TYPE_CHECKING, Optional
3 import numpy as np
4
5 import gymnasium as gym
6 from gymnasium import spaces, logger
7 from gymnasium.error import DependencyNotInstalled
8 from gymnasium.utils import EzPickle
9
10 from pettingzoo import AECEnv
11 from pettingzoo.utils import agent_selector
12
13 try:
14     import Box2D
15     from Box2D.b2 import (
16         circleShape,
17         contactListener,
18         edgeShape,
19         fixtureDef,
20         polygonShape,
21         revoluteJointDef,
22     )
23 except ImportError as e:
24     raise DependencyNotInstalled(
25         'Box2D is not installed, you can install it by run `pip install swig` '
26         + 'followed by `pip install "gymnasium[box2d]"`'
27     ) from e
28
29 if TYPE_CHECKING:
30     import pygame
31
32 FPS = 50
33 SCALE = 30.0 # affects how fast the game is, forces should be adjusted as well
34 MAIN_ENGINE_POWER = 13.0
35 SIDE_ENGINE_POWER = 0.6
36
37 INITIAL_RANDOM = 1000.0 # Set 1500 to make game harder
38
39 LANDER_POLY = [(-14, +17), (-17, 0), (-17, -10), (+17, -10), (+17, 0), (+14, +17)]
40 LEG_AWAY = 20
41 LEG_DOWN = 18
42 LEG_W, LEG_H = 2, 8
43 LEG_SPRING_TORQUE = 40
44
45 SIDE_ENGINE_HEIGHT = 14
46 SIDE_ENGINE_AWAY = 12
47 MAIN_ENGINE_Y_LOCATION = 4 # The Y loc of the main engine on the lander body.
48
49 VIEWPORT_W = 600
50 VIEWPORT_H = 400
51
52
53 class ContactDetector(contactListener):
54     def __init__(self, env):
55         contactListener.__init__(self)
56         self.env = env
57
58     def BeginContact(self, contact):
```

```

59     for lander in self.env.landings:
60         if (
61             lander == contact.fixtureA.body
62             or lander == contact.fixtureB.body
63         ):
64             self.env.game_over = True
65     for leg in self.env.legs:
66         if leg in [contact.fixtureA.body, contact.fixtureB.body]:
67             leg.ground_contact = True
68
69     def EndContact(self, contact):
70         for leg in self.env.legs:
71             if leg in [contact.fixtureA.body, contact.fixtureB.body]:
72                 leg.ground_contact = False
73
74
75 class SequentialEnv(AECEnv, EzPickle):
76     r"""
77     This is a rewrite of the Farama foundation Lunar Lander environment from:
78     https://gymnasium.farama.org/environments/box2d/lunar_lander/
79     Adapted for a AEC type petting zoo environment.
80     The action space and physics remain the same.
81     """
82     metadata = {
83         "render_modes": ["human", "rgb_array"],
84         "name": "multi_lander_v0",
85         "is_parallelizable": True,
86         "render_fps": FPS,
87     }
88
89     def __init__(
90         self,
91         render_mode: Optional[str] = None,
92         continuous: bool = False,
93         gravity: float = -10.0,
94         enable_wind: bool = False,
95         wind_power: float = 15.0,
96         turbulence_power: float = 1.5,
97         num_landers: int = 2,
98         *args, **kwargs
99     ):
100         EzPickle.__init__(
101             self,
102             render_mode,
103             continuous,
104             gravity,
105             enable_wind,
106             wind_power,
107             turbulence_power,
108             *args, **kwargs
109         )
110         AECEnv.__init__(self)
111         self.render_mode = render_mode
112         self.np_random = np.random
113
114         # Value Checkers
115         assert (-12.0 < gravity and gravity < 0.0
116             ), f"gravity (current value: {gravity}) must be between -12 and 0"
117         self.gravity = gravity
118

```

```

119 if 0.0 > wind_power or wind_power > 20.0:
120     logger.warn(
121         "wind_power value is recommended to be between 0.0 and 20.0, " +
122         f"(current value: {wind_power})"
123     )
124 self.wind_power = wind_power
125
126 if 0.0 > turbulence_power or turbulence_power > 2.0:
127     logger.warn(
128         "turbulence_power value is recommended to be between 0.0 and " +
129         f"2.0, (current value: {turbulence_power})"
130     )
131 # These are bounds for position realistically the environment
132 # should have ended long before we reach more than 50% outside
133 low = np.array(
134     [
135         -2.5, -2.5,                # x,y coordinates
136         -10.0, -10.0,             # x,y velocity bounds is 5x rated speed
137         -2 * math.pi, -10.0,     # Angle, Angular Velocity
138         -0.0, -0.0,              # L,R Leg not on ground
139     ]).astype(np.float32)
140 high = np.array(
141     [
142         2.5, 5.0, # 2.5,        # x,y coordinates
143         10.0, 10.0,             # x,y velocity bounds is 5x rated speed
144         2 * math.pi, 10.0,     # Angle, Angular Velocity
145         1.0, 1.0,               # L,R Leg on ground
146     ]).astype(np.float32)
147
148 # Environmental Variables
149 self.turbulence_power = turbulence_power
150 self.enable_wind = enable_wind
151 self.screen: pygame.Surface = None
152 self.clock = None
153 self.isopen = True
154 self.world = Box2D.b2World(gravity=(0, gravity))
155 self.moon = None
156 self.particles = []
157 self.landings = []
158 self.legs = []
159 self.prev_reward = None
160 self.continuous = continuous
161
162 # Agents
163 self.possible_agents = ["lander_" + str(i) for i in range(num_landings)]
164 self.agents = self.possible_agents[:]
165 self.agent_name_mapping = dict(zip(self.agents,
166                                     list(range(self.num_agents))))
167 self._agent_selector = agent_selector(self.agents)
168
169 # Spaces
170 self.observation_spaces = dict(zip(self.agents,
171                                     [spaces.Box(low, high)] * self.num_agents))
172 if self.continuous:
173     # Action is two floats [main engine, left-right engines].
174     # Main engine: -1..0 off, 0..+1 throttle from 50% to 100% power.
175     # Engine can't work with less than 50% power.
176     # Left-right: -1.0..-0.5 fire left engine, +0.5..+1.0 fire right
177     # engine, -0.5..0.5 off
178     self._action_space = spaces.Box(-1, +1, (2,), dtype=np.float32)

```

```

179     else:
180         # No-op, fire left engine, main engine, right engine
181         self._action_space = spaces.Discrete(4)
182     self.action_spaces = dict(zip(self.agents,
183                                   [self._action_space]*self.num_agents))

```

```

185 def observation_space(self, agent):
186     return self.observation_spaces[agent]

```

```

188 def action_space(self, agent):
189     return self.action_spaces[agent]

```

```

191 def convert_to_dict(self, list_of_list):
192     return dict(zip(self.agents, list_of_list))

```

```

194 def close(self):
195     if self.screen is not None:
196         import pygame
197         pygame.display.quit()
198         pygame.quit()
199         pygame.QUIT
200     self.isopen = False

```

```

202 def create_landers(self):
203     landers = []
204     for n in range(self.num_agents):
205         lander: Box2D.b2Body = self.world.CreateDynamicBody(
206             position=((n+1)/(self.num_agents+1) * VIEWPORT_W / SCALE,
207                       VIEWPORT_H / SCALE),
208             angle=0.0,
209             fixtures=fixtureDef(
210                 shape=polygonShape(
211                     vertices=[(x/SCALE, y/SCALE) for x, y in LANDER_POLY]),
212                 density=5.0,
213                 friction=0.1,
214                 categoryBits=0x0010,
215                 #maskBits=0x001, # uncomment for : collide only with ground
216                 restitution=0.0,
217             ), # 0.99 bouncy
218         )
219         lander.color1 = ((128), (102+75*(n))%255, (230))
220         lander.color2 = ((77), (77), (128))
221         landers.append(lander)
222     self.landerns = landers

```

```

224 legs = []
225 for n, lander in enumerate(self.landerns):
226     for i in [-1, +1]:
227         leg = self.world.CreateDynamicBody(
228             position=(lander.position[0] - i * LEG_AWAY / SCALE,
229                     lander.position[1]),
230             angle=(i * 0.05),
231             fixtures=fixtureDef(
232                 shape=polygonShape(box=(LEG_W / SCALE, LEG_H / SCALE)),
233                 density=1.0,
234                 restitution=0.0,
235                 categoryBits=0x0020,
236                 #maskBits=0x001,
237             ),
238         )

```

```

239         leg.ground_contact = False
240         leg.color1 = ((128), (102+75*(n))%255, (230))
241         leg.color2 = ((77), (77), (128))
242         rjd = revoluteJointDef(
243             bodyA=lander,
244             bodyB=leg,
245             localAnchorA=(0, 0),
246             localAnchorB=(i * LEG_AWAY / SCALE, LEG_DOWN / SCALE),
247             enableMotor=True,
248             enableLimit=True,
249             maxMotorTorque=LEG_SPRING_TORQUE,
250             motorSpeed=+0.3 * i, # low enough not to jump into the sky
251         )
252         if i == -1:
253             rjd.lowerAngle = (
254                 +0.9 - 0.5
255             ) # Valid angle of travel for the legs
256             rjd.upperAngle = +0.9
257         else:
258             rjd.lowerAngle = -0.9
259             rjd.upperAngle = -0.9 + 0.5
260         leg.joint = self.world.CreateJoint(rjd)
261         legs.append(leg)
262     self.legs = legs
263
264     def _create_particle(self, mass, x, y, ttl):
265         p = self.world.CreateDynamicBody(
266             position=(x, y),
267             angle=0.0,
268             fixtures=fixtureDef(
269                 shape=circleShape(radius=2 / SCALE, pos=(0, 0)),
270                 density=mass,
271                 friction=0.1,
272                 categoryBits=0x0100,
273                 maskBits=0x001, # collide only with ground
274                 restitution=0.3,
275             ),
276         )
277         p.ttl = ttl
278         self.particles.append(p)
279         self._clean_particles(False)
280         return p
281
282     def _clean_particles(self, all_particle):
283         while self.particles and (all_particle or self.particles[0].ttl < 0):
284             self.world.DestroyBody(self.particles.pop(0))
285
286
287     def render(self):
288         if self.render_mode is None:
289             assert self.spec is not None
290             gym.logger.warn(
291                 "You are calling render method without specifying render mode."
292                 "You can specify the render_mode at initialization, "
293                 f'e.g. gym.make("{self.spec.id}", render_mode="rgb_array")')
294             return
295
296         try:
297             import pygame
298             from pygame import gfxdraw

```



```

299 except ImportError as e:
300     raise DependencyNotInstalled(
301         'pygame is not installed, run `pip install "gymnasium[box2d]"`
302     ) from e
303
304 if self.screen is None and self.render_mode == "human":
305     pygame.init()
306     pygame.display.init()
307     self.screen = pygame.display.set_mode((VIEWPORT_W, VIEWPORT_H))
308 if self.clock is None:
309     self.clock = pygame.time.Clock()
310
311 self.surf = pygame.Surface((VIEWPORT_W, VIEWPORT_H))
312 pygame.transform.scale(self.surf, (SCALE, SCALE))
313 pygame.draw.rect(self.surf, (255, 255, 255), self.surf.get_rect())
314
315 for obj in self.particles:
316     obj.ttl -= 0.15
317     obj.color1 = (
318         int(max(0.2, 0.15 + obj.ttl) * 255),
319         int(max(0.2, 0.5 * obj.ttl) * 255),
320         int(max(0.2, 0.5 * obj.ttl) * 255),
321     )
322     obj.color2 = (
323         int(max(0.2, 0.15 + obj.ttl) * 255),
324         int(max(0.2, 0.5 * obj.ttl) * 255),
325         int(max(0.2, 0.5 * obj.ttl) * 255),
326     )
327
328 self._clean_particles(False)
329
330 for p in self.sky_polys:
331     scaled_poly = []
332     for coord in p:
333         scaled_poly.append((coord[0] * SCALE, coord[1] * SCALE))
334     pygame.draw.polygon(self.surf, (0, 0, 0), scaled_poly)
335     gfxdraw.aapolygon(self.surf, scaled_poly, (0, 0, 0))
336
337 for obj in self.particles + self.drawlist:
338     for f in obj.fixtures:
339         trans = f.body.transform
340         if type(f.shape) is circleShape:
341             pygame.draw.circle(
342                 self.surf,
343                 color=obj.color1,
344                 center=trans * f.shape.pos * SCALE,
345                 radius=f.shape.radius * SCALE,
346             )
347             pygame.draw.circle(
348                 self.surf,
349                 color=obj.color2,
350                 center=trans * f.shape.pos * SCALE,
351                 radius=f.shape.radius * SCALE,
352             )
353         else:
354             path = [trans * v * SCALE for v in f.shape.vertices]
355             pygame.draw.polygon(self.surf, color=obj.color1, points=path)
356             gfxdraw.aapolygon(self.surf, path, obj.color1)
357             pygame.draw.aalines(
358                 self.surf, color=obj.color2, points=path, closed=True)

```

359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418

```
for x in [self.helipad_x1, self.helipad_x2]:
    x = x * SCALE
    flagy1 = self.helipad_y * SCALE
    flagy2 = flagy1 + 50
    pygame.draw.line(
        self.surf,
        color=(255, 255, 255),
        start_pos=(x, flagy1),
        end_pos=(x, flagy2),
        width=1,
    )
    pygame.draw.polygon(
        self.surf,
        color=(204, 204, 0),
        points=[
            (x, flagy2),
            (x, flagy2 - 10),
            (x + 25, flagy2 - 5),
        ],
    )
    gfxdraw.aapolygon(
        self.surf,
        [(x, flagy2), (x, flagy2 - 10), (x + 25, flagy2 - 5)],
        (204, 204, 0),
    )
```

```
self.surf = pygame.transform.flip(self.surf, False, True)
```

```
if self.render_mode == "human":
    assert self.screen is not None
    self.screen.blit(self.surf, (0, 0))
    pygame.event.pump()
    self.clock.tick(self.metadata["render_fps"])
    pygame.display.flip()
elif self.render_mode == "rgb_array":
    return np.transpose(
        np.array(pygame.surfarray.pixels3d(self.surf)), axes=(1, 0, 2)
    )
```

```
# End Render()
```

```
def _destroy(self):
    if not self.moon:
        return
    self.world.contactListener = None
    self._clean_particles(True)
    self.world.DestroyBody(self.moon)
    self.moon = None
    for lander in self.landings:
        self.world.DestroyBody(lander)
        lander = None
    for i in range(len(self.legs)):
        self.world.DestroyBody(self.legs[i])
```

```
def reset(self, *, seed: Optional[int] = None, options: Optional[dict] = None):
    self._destroy()
    # Issue: https://github.com/Farama-Foundation/Gymnasium/issues/728
    # self._destroy() is not enough to clean(reset), workaround is
    # to create a totally new world for self.reset()
```

```

419 self.world = Box2D.b2World(gravity=(0, self.gravity))
420 self.world.contactListener_keepref = ContactDetector(self)
421 self.world.contactListener = self.world.contactListener_keepref
422 self.game_over = False
423 self.prev_shaping = [None]*self.num_agents
424
425 W = VIEWPORT_W / SCALE
426 H = VIEWPORT_H / SCALE
427
428 # Create Terrain
429 CHUNKS = 11
430 height = self.np_random.uniform(0, H / 2, size=(CHUNKS + 1,))
431 chunk_x = [W / (CHUNKS - 1) * i for i in range(CHUNKS)]
432 self.helipad_x1 = chunk_x[CHUNKS // 2 - 1]
433 self.helipad_x2 = chunk_x[CHUNKS // 2 + 1]
434 self.helipad_y = H / 4
435 height[CHUNKS // 2 - 2] = self.helipad_y
436 height[CHUNKS // 2 - 1] = self.helipad_y
437 height[CHUNKS // 2 + 0] = self.helipad_y
438 height[CHUNKS // 2 + 1] = self.helipad_y
439 height[CHUNKS // 2 + 2] = self.helipad_y
440 smooth_y = [
441     0.33 * (height[i - 1] + height[i + 0] + height[i + 1])
442     for i in range(CHUNKS)
443 ]
444
445 self.moon = self.world.CreateStaticBody(
446     shapes=edgeShape(vertices=[(0, 0), (W, 0)])
447 )
448 self.sky_polys = []
449 for i in range(CHUNKS - 1):
450     p1 = (chunk_x[i], smooth_y[i])
451     p2 = (chunk_x[i + 1], smooth_y[i + 1])
452     self.moon.CreateEdgeFixture(vertices=[p1,p2],density=0,friction=0.1)
453     self.sky_polys.append([p1, p2, (p2[0], H), (p1[0], H)])
454
455 self.moon.color1 = (0.0, 0.0, 0.0)
456 self.moon.color2 = (0.0, 0.0, 0.0)
457
458 if self.enable_wind: # Initialize wind pattern based on index
459     self.wind_idx = self.np_random.integers(-9999, 9999)
460     self.torque_idx = self.np_random.integers(-9999, 9999)
461
462 self.create_landers()
463
464 for lander in self.landings:
465     lander.ApplyForceToCenter(
466         (
467             self.np_random.uniform(-INITIAL_RANDOM, INITIAL_RANDOM),
468             self.np_random.uniform(-INITIAL_RANDOM, INITIAL_RANDOM),
469         ),
470         True,
471     )
472 self.drawlist = self.legs + self.landings
473 if self.render_mode == "human": self.render()
474
475 self.agents = self.possible_agents[:]
476 self._agent_selector = agent_selector(self.agents)
477 self.agent_selection = self._agent_selector.next()
478 self.rewards = {agent: 0 for agent in self.agents}

```

```

479     self._cumulative_rewards = {agent: 0 for agent in self.agents}
480     self.terminations = {agent: False for agent in self.agents}
481     self.truncations = {agent: False for agent in self.agents}
482     self.infos = {agent: {} for agent in self.agents}
483
484     return self.observe(), None
485 # End reset()
486
487
488 def observe(self, agent=None):
489     if not agent : agent = self.agent_selection
490     id = self.agent_name_mapping[agent]
491     lander = self.landings[id]
492     pos = lander.position
493     vel = lander.linearVelocity
494
495     observation = [
496         (pos.x - VIEWPORT_W/SCALE/2) / (VIEWPORT_W/SCALE/2),
497         (pos.y - (self.helipad_y + LEG_DOWN/SCALE)) / (VIEWPORT_H/SCALE/2),
498         vel.x * (VIEWPORT_W/SCALE/2) / FPS,
499         vel.y * (VIEWPORT_H/SCALE/2) / FPS,
500         lander.angle,
501         20.0 * lander.angularVelocity / FPS,
502         1.0 if self.legs[id*2+0].ground_contact else 0.0,
503         1.0 if self.legs[id*2+1].ground_contact else 0.0,
504     ]
505     assert len(observation) == 8
506     return observation
507
508 def latest_reward_state(self, agent, state):
509     reward = 0
510     id = self.agent_name_mapping[agent]
511     shaping = (
512         -100 * np.sqrt(state[0] * state[0] + state[1] * state[1])
513         - 100 * np.sqrt(state[2] * state[2] + state[3] * state[3])
514         - 100 * abs(state[4])
515         + 10 * state[6] + 10 * state[7]
516     ) # And ten points for legs contact, the idea is if you
517        # lose contact again after landing, you get negative reward
518     if self.prev_shaping[id] is not None:
519         reward = shaping - self.prev_shaping[id]
520     self.prev_shaping[id] = shaping
521     return reward
522
523
524 def step(self, action):
525     assert self.landings is not None, "You forgot to call reset()"
526     if (
527         self.terminations[self.agent_selection]
528         or self.truncations[self.agent_selection]
529     ):
530         # If one agent has terminated this accepts a None action,
531         # which otherwise errors, handles stepping to the next agent
532         obs = self.observe(self.agent_selection)
533         self.agent_selection = self._agent_selector.next()
534         return obs, 0, True, False, {}
535
536     # Update wind and apply to the lander
537     if self.enable_wind and not any([l.ground_contact for l in self.legs]):
538         # the function used for wind is tanh(sin(2 k x) + sin(pi k x)),

```

```

539     # which is proven to never be periodic, k = 0.01
540     wind_mag = (
541         math.tanh(
542             math.sin(0.02 * self.wind_idx)
543             + (math.sin(math.pi * 0.01 * self.wind_idx))
544         ) * self.wind_power
545     )
546     self.wind_idx += 1
547
548     for lander in self.landings:
549         lander.ApplyForceToCenter(
550             (wind_mag, 0.0),
551             True,
552         )
553
554     # the function used for torque is tanh(sin(2 k x) + sin(pi k x)),
555     # which is proven to never be periodic, k = 0.01
556     torque_mag = (
557         math.tanh(
558             math.sin(0.02 * self.torque_idx)
559             + (math.sin(math.pi * 0.01 * self.torque_idx))
560         ) * self.turbulence_power
561     )
562     self.torque_idx += 1
563
564     for lander in self.landings:
565         lander.ApplyTorque(
566             torque_mag,
567             True,
568         )
569
570     # For current agent:
571     agent = self.agent_selection
572     lander = self.landings[self.agent_name_mapping[agent]]
573     # Check action validity
574     if self.continuous:
575         action = np.clip(action, -1, +1).astype(np.float64)
576     else:
577         assert self.action_spaces[agent].contains(
578             action), f"{action!r} ({type(action)}) invalid "
579
580     # Tip is the (X and Y) components of the rotation of the lander.
581     tip = (math.sin(lander.angle), math.cos(lander.angle))
582
583     # Side is the (-Y and X) components of the rotation of the lander.
584     side = (-tip[1], tip[0])
585
586     # Generate two random numbers between -1/SCALE and 1/SCALE.
587     dispersion = [self.np_random.uniform(-1.0,+1.0)/SCALE for _ in range(2)]
588
589     m_power = 0.0
590     if (self.continuous and action[0] > 0.0) or (
591         not self.continuous and action == 2
592     ):
593         # Main engine
594         if self.continuous:
595             m_power = (np.clip(action[0], 0.0, 1.0) + 1.0) * 0.5 # 0.5..1.0
596             assert m_power >= 0.5 and m_power <= 1.0
597         else:
598             m_power = 1.0

```

```

599
600     # 4 is move a bit downwards, +-2 for randomness
601     # The components of the impulse to be applied by the main engine.
602     ox = (
603         tip[0] * (MAIN_ENGINE_Y_LOCATION / SCALE + 2 * dispersion[0])
604         + side[0] * dispersion[1]
605     )
606     oy = (
607         -tip[1] * (MAIN_ENGINE_Y_LOCATION / SCALE + 2 * dispersion[0])
608         - side[1] * dispersion[1]
609     )
610     impulse_pos = (lander.position[0] + ox, lander.position[1] + oy)
611     if self.render_mode is not None:
612         # particles are just a decoration, with no physics impact,
613         # so don't add them when not rendering
614         p = self._create_particle(
615             3.5, # 3.5 is here to make particle speed adequate
616             impulse_pos[0],
617             impulse_pos[1],
618             m_power,
619         )
620         p.ApplyLinearImpulse(
621             (
622                 ox * MAIN_ENGINE_POWER * m_power,
623                 oy * MAIN_ENGINE_POWER * m_power,
624             ),
625             impulse_pos,
626             True,
627         )
628         lander.ApplyLinearImpulse(
629             (-ox * MAIN_ENGINE_POWER * m_power,
630              -oy * MAIN_ENGINE_POWER * m_power),
631             impulse_pos,
632             True,
633         )
634
635     s_power = 0.0
636     if (self.continuous and np.abs(action[1]) > 0.5) or (
637         not self.continuous and action in [1, 3]
638     ):
639         # Orientation/Side engines
640         if self.continuous:
641             direction = np.sign(action[1])
642             s_power = np.clip(np.abs(action[1]), 0.5, 1.0)
643             assert s_power >= 0.5 and s_power <= 1.0
644         else:
645             # action = 1 is left, action = 3 is right
646             direction = action - 2
647             s_power = 1.0
648
649     # The components of the impulse to be applied by the side engines.
650     ox = tip[0] * dispersion[0] + side[0] * (
651         3 * dispersion[1] + direction * SIDE_ENGINE_AWAY / SCALE
652     )
653     oy = -tip[1] * dispersion[0] - side[1] * (
654         3 * dispersion[1] + direction * SIDE_ENGINE_AWAY / SCALE
655     )
656
657     # The constant 17 is presumably meant to be SIDE_ENGINE_HEIGHT.
658     # However, SIDE_ENGINE_HEIGHT is defined as 14, causing the

```

```

659     # position of the thrust on the body of the lander to change,
660     # depending on the orientation of the lander. This results in
661     # an orientation dependent torque being applied to the lander.
662
663     impulse_pos = (
664         lander.position[0] + ox - tip[0] * 17 / SCALE,
665         lander.position[1] + oy + tip[1] * SIDE_ENGINE_HEIGHT / SCALE,
666     )
667     if self.render_mode is not None:
668         # particles are just decoration, with no impact on the physics,
669         # so don't add them when not rendering
670         p = self._create_particle(0.7, impulse_pos[0],
671                                   impulse_pos[1], s_power)
672         p.ApplyLinearImpulse(
673             (ox * SIDE_ENGINE_POWER * s_power,
674              oy * SIDE_ENGINE_POWER * s_power,),
675             impulse_pos,
676             True,
677         )
678     lander.ApplyLinearImpulse(
679         (-ox * SIDE_ENGINE_POWER * s_power,
680          -oy * SIDE_ENGINE_POWER * s_power),
681         impulse_pos,
682         True,
683     )
684
685     # Update Positions
686     self.world.Step(1.0 / FPS, 6 * 30, 2 * 30)
687
688     observation = self.observe(agent)
689     reward = self.latest_reward_state(agent, observation)
690     reward -= m_power * 0.30 # Deduct cost of fuel
691     reward -= s_power * 0.03 # Deduct cost of fuel
692
693     if self.game_over or abs(observation[0]) >= 1.0:
694         self.terminations[self.agent_selection] = True
695         reward = -100
696     if not lander.awake:
697         self.terminations[self.agent_selection] = True
698         reward = +100
699     terminated = self.terminations[self.agent_selection]
700
701     self.rewards[agent] = reward
702     self._accumulate_rewards()
703
704     self.agent_selection = self._agent_selector.next()
705
706     if self.render_mode == "human":
707         self.render()
708     # truncation=False as the time limit is handled by the `TimeLimit`
709     return observation, reward, terminated, False, {}
710
711
712 class Parallel_Env(SequentialEnv):
713     def observe(self):
714         obs_list = []
715         for i, lander in enumerate(self.landings):
716             pos = lander.position
717             vel = lander.linearVelocity
718             obs = [

```

```

719         (pos.x - VIEWPORT_W/SCALE/2) / (VIEWPORT_W/SCALE/2),
720         (pos.y - (self.helipad_y + LEG_DOWN/SCALE)) / (VIEWPORT_H/SCALE/2),
721         vel.x * (VIEWPORT_W/SCALE/2) / FPS,
722         vel.y * (VIEWPORT_H/SCALE/2) / FPS,
723         lander.angle,
724         20.0 * lander.angularVelocity / FPS,
725         1.0 if self.legs[i*2+0].ground_contact else 0.0,
726         1.0 if self.legs[i*2+1].ground_contact else 0.0,
727     ]
728     assert len(obs) == 8
729     obs_list.append(obs)
730     return dict(zip(self.possible_agents, obs_list))
731
732 def last(self):
733     return self.observe()
734
735 def step(self, action_list):
736     self._clear_rewards()
737     assert self.landings is not None, "You forgot to call reset()"
738
739     # Update wind and apply to the lander
740     if self.enable_wind and not any([l.ground_contact for l in self.legs]):
741         # the function used for wind is  $\tanh(\sin(2 k x) + \sin(\pi k x))$ ,
742         # which is proven to never be periodic,  $k = 0.01$ 
743         wind_mag = (
744             math.tanh(
745                 math.sin(0.02 * self.wind_idx)
746                 + (math.sin(math.pi * 0.01 * self.wind_idx))
747             ) * self.wind_power
748         )
749         self.wind_idx += 1
750
751         for lander in self.landings:
752             lander.ApplyForceToCenter(
753                 (wind_mag, 0.0),
754                 True,
755             )
756         # the function used for torque is  $\tanh(\sin(2 k x) + \sin(\pi k x))$ ,
757         # which is proven to never be periodic,  $k = 0.01$ 
758         torque_mag = (
759             math.tanh(
760                 math.sin(0.02 * self.torque_idx)
761                 + (math.sin(math.pi * 0.01 * self.torque_idx))
762             ) * self.turbulence_power
763         )
764         self.torque_idx += 1
765
766         for lander in self.landings:
767             lander.ApplyTorque(
768                 torque_mag,
769                 True,
770             )
771
772     #for agent in self.agents:
773     for agent, action in action_list.items():
774         # For current agent:
775         lander = self.landings[self.agent_name_mapping[agent]]
776         # Check action validity
777         if self.continuous:
778             action = np.clip(action, -1, +1).astype(np.float64)

```



```

779 else:
780     assert self.action_spaces[agent].contains(
781         action), f"{action!r} ({type(action)}) invalid "
782
783     # Tip is the (X and Y) components of the rotation of the lander.
784     tip = (math.sin(lander.angle), math.cos(lander.angle))
785
786     # Side is the (-Y and X) components of the rotation of the lander.
787     side = (-tip[1], tip[0])
788
789     # Generate two random numbers between -1/SCALE and 1/SCALE.
790     dispersion = [self.np_random.uniform(-1.0,+1.0)/SCALE for _ in range(2)]
791
792     m_power = 0.0
793     if (self.continuous and action[0] > 0.0) or (
794         not self.continuous and action == 2
795     ):
796         # Main engine
797         if self.continuous:
798             m_power = (np.clip(action[0], 0.0, 1.0) + 1.0) * 0.5 # 0.5..1.0
799             assert m_power >= 0.5 and m_power <= 1.0
800         else:
801             m_power = 1.0
802
803         # 4 is move a bit downwards, +-2 for randomness
804         # The components of the impulse to be applied by the main engine.
805         ox = (
806             tip[0] * (MAIN_ENGINE_Y_LOCATION / SCALE + 2 * dispersion[0])
807             + side[0] * dispersion[1]
808         )
809         oy = (
810             -tip[1] * (MAIN_ENGINE_Y_LOCATION / SCALE + 2 * dispersion[0])
811             - side[1] * dispersion[1]
812         )
813         impulse_pos = (lander.position[0] + ox, lander.position[1] + oy)
814         if self.render_mode is not None:
815             # particles are just a decoration, with no physics impact,
816             # so don't add them when not rendering
817             p = self._create_particle(
818                 3.5, # 3.5 is here to make particle speed adequate
819                 impulse_pos[0],
820                 impulse_pos[1],
821                 m_power,
822             )
823             p.ApplyLinearImpulse(
824                 (
825                     ox * MAIN_ENGINE_POWER * m_power,
826                     oy * MAIN_ENGINE_POWER * m_power,
827                 ),
828                 impulse_pos,
829                 True,
830             )
831             lander.ApplyLinearImpulse(
832                 (-ox * MAIN_ENGINE_POWER * m_power,
833                  -oy * MAIN_ENGINE_POWER * m_power),
834                 impulse_pos,
835                 True,
836             )
837
838     s_power = 0.0

```

```

839     if (self.continuous and np.abs(action[1]) > 0.5) or (
840         not self.continuous and action in [1, 3]
841     ):
842         # Orientation/Side engines
843         if self.continuous:
844             direction = np.sign(action[1])
845             s_power = np.clip(np.abs(action[1]), 0.5, 1.0)
846             assert s_power >= 0.5 and s_power <= 1.0
847         else:
848             # action = 1 is left, action = 3 is right
849             direction = action - 2
850             s_power = 1.0
851
852         # The components of the impulse to be applied by the side engines.
853         ox = tip[0] * dispersion[0] + side[0] * (
854             3 * dispersion[1] + direction * SIDE_ENGINE_AWAY / SCALE
855         )
856         oy = -tip[1] * dispersion[0] - side[1] * (
857             3 * dispersion[1] + direction * SIDE_ENGINE_AWAY / SCALE
858         )
859
860         # The constant 17 is presumably meant to be SIDE_ENGINE_HEIGHT.
861         # However, SIDE_ENGINE_HEIGHT is defined as 14, causing the
862         # position of the thrust on the body of the lander to change,
863         # depending on the orientation of the lander. This results in
864         # an orientation dependent torque being applied to the lander.
865
866         impulse_pos = (
867             lander.position[0] + ox - tip[0] * 17 / SCALE,
868             lander.position[1] + oy + tip[1] * SIDE_ENGINE_HEIGHT / SCALE,
869         )
870         if self.render_mode is not None:
871             # particles are just decoration, with no impact on the physics,
872             # so don't add them when not rendering
873             p = self._create_particle(0.7, impulse_pos[0],
874                                     impulse_pos[1], s_power)
875             p.ApplyLinearImpulse(
876                 (ox * SIDE_ENGINE_POWER * s_power,
877                  oy * SIDE_ENGINE_POWER * s_power),
878                 impulse_pos,
879                 True,
880             )
881             lander.ApplyLinearImpulse(
882                 (-ox * SIDE_ENGINE_POWER * s_power,
883                  -oy * SIDE_ENGINE_POWER * s_power),
884                 impulse_pos,
885                 True,
886             )
887             self.rewards[agent] -= m_power * 0.30 # Deduct cost of fuel
888             self.rewards[agent] -= s_power * 0.03 # Deduct cost of fuel
889 """ End Agent Loop """
890
891 # Update Positions
892 self.world.Step(1.0 / FPS, 6 * 30, 2 * 30)
893
894 obs_list = self.observe()
895
896 for agent in self.agents:
897     obs = obs_list[agent]
898     lander = self.landings[self.agent_name_mapping[agent]]

```

```

899         self.rewards[agent] += self.latest_reward_state(agent, obs)
900         if self.game_over or abs(obs[0]) >= 1.0:
901             self.terminations[agent] = True
902             self.rewards[agent] -= 100
903         if not lander.awake:
904             self.terminations[agent] = True
905             self.rewards[agent] += 100
906
907     self._accumulate_rewards()
908
909     if self.render_mode == "human":
910         self.render()
911         # truncation=False as the time limit is handled by the `TimeLimit`
912     return obs_list, self.rewards, self.terminations, self.truncations, {}
913
914
915 def heuristic(env, s):
916     """
917     The heuristic for
918     1. Testing
919     2. Demonstration rollout.
920
921     Args:
922         env: The environment
923         s (list): The state. Attributes:
924             s[0] is the horizontal coordinate
925             s[1] is the vertical coordinate
926             s[2] is the horizontal speed
927             s[3] is the vertical speed
928             s[4] is the angle
929             s[5] is the angular speed
930             s[6] 1 if first leg has contact, else 0
931             s[7] 1 if second leg has contact, else 0
932
933     Returns:
934         a: The heuristic to be fed into the step function defined above to
935            determine the next step and reward.
936     """
937
938     angle_targ = s[0] * 0.5 + s[2] * 1.0 # angle should point towards center
939     if angle_targ > 0.4:
940         angle_targ = 0.4 # more than 0.4 radians (22 degrees) is bad
941     if angle_targ < -0.4:
942         angle_targ = -0.4
943     hover_targ = 0.55 * np.abs(
944         s[0]
945     ) # target y should be proportional to horizontal offset
946
947     angle_todo = (angle_targ - s[4]) * 0.5 - (s[5]) * 1.0
948     hover_todo = (hover_targ - s[1]) * 0.5 - (s[3]) * 0.5
949
950     if s[6] or s[7]: # legs have contact
951         angle_todo = 0
952         hover_todo = (
953             -(s[3]) * 0.5
954         ) # override to reduce fall speed, that's all we need after contact
955
956     if env.unwrapped.continuous:
957         a = np.array([hover_todo * 20 - 1, -angle_todo * 20])
958         a = np.clip(a, -1, +1)

```

```

959     else:
960         a = 0
961         if hover_todo > np.abs(angle_todo) and hover_todo > 0.05:
962             a = 2
963         elif angle_todo < -0.05:
964             a = 3
965         elif angle_todo > +0.05:
966             a = 1
967     return a
968
969 def demo_heuristic_lander(env, reps=1, seed=None, render=False):
970     total_reward = 0
971     steps = 0
972     for _ in range(reps):
973         s = env.reset(seed=seed)
974         while True:
975             s = env.last()[0]
976             a = heuristic(env, s)
977             s, r, terminated, truncated, info = env.step(a)
978             total_reward += r
979
980             if render:
981                 still_open = env.render()
982                 if still_open is False:
983                     break
984
985             if steps % 20 == 0 or terminated or truncated:
986                 print("observations:", " ".join([f"{x:+0.2f}" for x in s]))
987                 print(f"step {steps} total_reward {total_reward:+0.2f}")
988                 print(env.terminations.values())
989             steps += 1
990             # if terminated or truncated:
991             if all(env.terminations.values()):
992                 break
993     if render:
994         env.close()
995     return total_reward
996
997 def demo_parallel_heuristic(env, reps=1, seed=None, render=False):
998     total_reward = 0
999     steps = 0
1000    for _ in range(reps):
1001        _ = env.reset(seed=seed)
1002        obs = env.last()
1003        while True:
1004            actions = {agent: heuristic(env, obs[agent]) for agent in env.agents}
1005            obs, r, terminateds, truncateds, info = env.step(actions)
1006            total_reward += sum(r.values())
1007
1008            if render:
1009                still_open = env.render()
1010                if still_open is False:
1011                    break
1012
1013            if steps % 20 == 0 or all(terminateds) or all(truncateds):
1014                for agent in env.agents:
1015                    print(f"{agent} observations:", " ".join([f"{x:+0.2f}" for x in obs[ag
1016                    print(f"step {steps} total_reward {total_reward:+0.2f}")
1017                    print(env.terminations.values())
1018            steps += 1

```

```

1019         # if terminated or truncated:
1020         if all(env.terminations.values()):
1021             break
1022     if render:
1023         env.close()
1024     return total_reward
1025
1026 def parallel_env(**kwargs):
1027     return Parallel_Env(**kwargs)
1028
1029 def sequential_env(**kwargs):
1030     return SequentialEnv(**kwargs)
1031
1032 def env(**kwargs):
1033     return Parallel_Env(**kwargs)
1034
1035 if __name__ == "__main__":
1036     import argparse
1037     parser = argparse.ArgumentParser(description=
1038                                     'A multi-agent version of Lunar Lander')
1039     parser.add_argument('-s', '--sequential', action="store_true",
1040                        help='Iterate as AEC environment')
1041     parser.add_argument('-n', '--num_landers', type=int, default=2,
1042                        help='number of landers')
1043     parser.add_argument('-d', '--demo_iters', type=int, default=1,
1044                        help='number of landers')
1045     args = parser.parse_args()
1046
1047     if args.sequential:
1048         _env = sequential_env(render_mode="human", num_landers=args.num_landers)
1049         demo_heuristic_lander(_env, reps=args.demo_iters, render=True)
1050     else:
1051         _env = parallel_env(render_mode="human", num_landers=args.num_landers)
1052         demo_parallel_heuristic(_env, reps=args.demo_iters, render=True)

```