# Welcome to:

# Hands-on RL with Ray's RLlib

A beginner's tutorial for working with environments, models, and algorithms

# Who am I?

Sven Mika - ML Engineer Anyscale Inc.



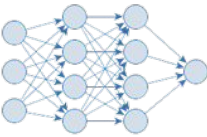https://linkedin.com/in/sven-mika
sven@anyscale.com

- Late 90s: Biochemistry (undergrad)
- Early-mid 2000s: Bioinformatics & NLP (PhD)

  *"Why are you using NNs? SVMs are so much better and robust!"*

- 2010s: Wall Street (Quant. Dev./Data Scientist)
- 2015+: Contributor to other OSS RL libraries (TensorForce, RLgraph, Surreal); Self-taught RL: Books, papers, online courses.
- 2019+: Anyscale Inc. (lead-dev RLlib)

**Setup:** conda create -n rllib python=3.8; conda activate rllib; pip install ray[rllib]==1.4; pip install [tensorflow|torch]; pip install jupyter-labs; git clone https://github.com/sven1977/rllib_tutorials; cd rllib_tutorials; jupyter-lab
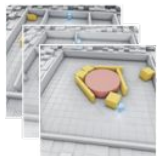
policy (in deep-RL, this is a neural network)

tf.keras / torch.nn.Modules

multi-GPU

agents

built-in LSTM + attention nets

observations

arbitrary obs spaces

15+ available algorithms (model-free; model-based; offline RL)

actions

multi-agent

rewards

+0.3
-1.0

reward shaping / learning w/o rewards

external environments

arbitrary act. spaces

custom callbacks

environment

parallelize and distribute

# Key differences: RL vs Supervised Learning

- Data collection loop is essential part of the RL algorithm, especially in the distributed setting.

- RL is forced to use unstable (bootstrapped) loss functions due to missing labels (rewards cannot cover absence of SL-style labels).

$$\mathbb{E}_{(s,a,r,s')}\left[\frac{1}{2}\left(\underbrace{R(s,a,s') + \gamma\max_{a'}Q(s',a';\theta)}_{\text{``labels'' (bootstrapped using reward AND network output)}} - \underbrace{Q(s,a;\theta)}_{\text{predictions (actual network output)}}\right)^2\right]$$

- Exploration vs exploitation. An (online) RL algo can try new things in the environment.

- By default, RL tries to optimize over a time axis.
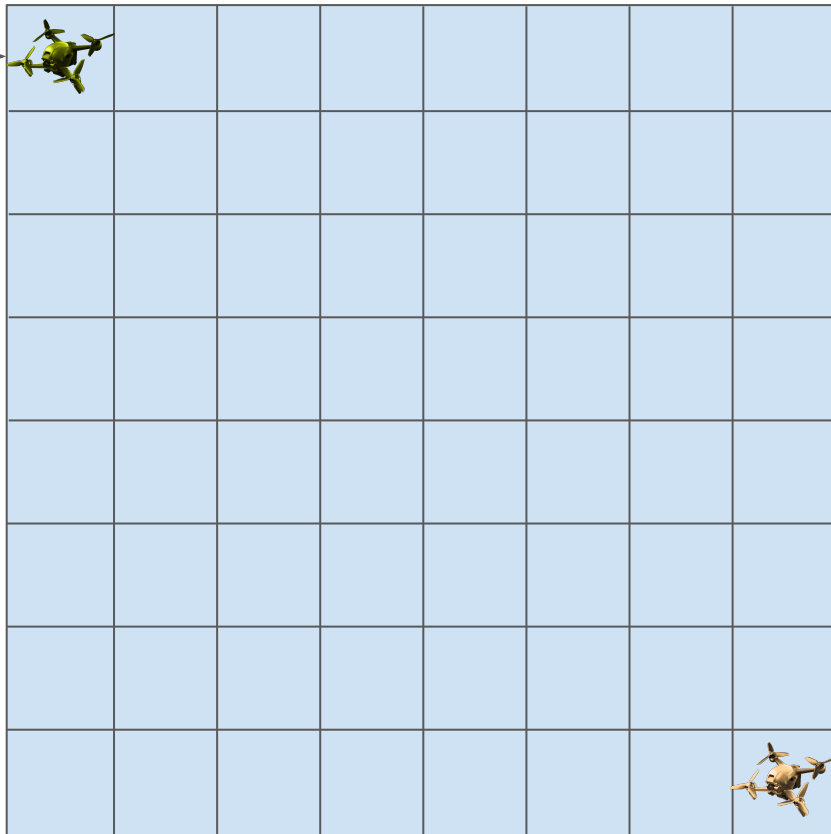
# Overview of RLlib's Industry Users

# Many Great Reasons for Using RLlib

-   RLlib is based on Ray, so it benefits from all its improvements on **performance** and **scalability**.

-   RLlib is backed by Anyscale, the fast-growing, well-funded company behind Ray, offering strong **OSS support**.

-   RLlib is extremely **flexible**, allowing you to **customize every aspect** of the RL cycle and workflows.

-   RLlib is good at solving **real-world** problems, supporting **offline RL**, **multi-agent** setups, **external simulators**, and more.

# Our Environment - The "MultiAgentArena" 🙂


agent 1

action space: Discrete(4)
0 - up
1 - right
2 - down
3 - left

rewards:
agent1 ("cover as much ground as possible"):
● +1 if new field
● -1 if colliding with agent 2
● -0.5 otherwise

agent2 ("defend: bump into agent1 as often as possible"):
● +1 if colliding with agent 1
● -0.1 otherwise

observation space:
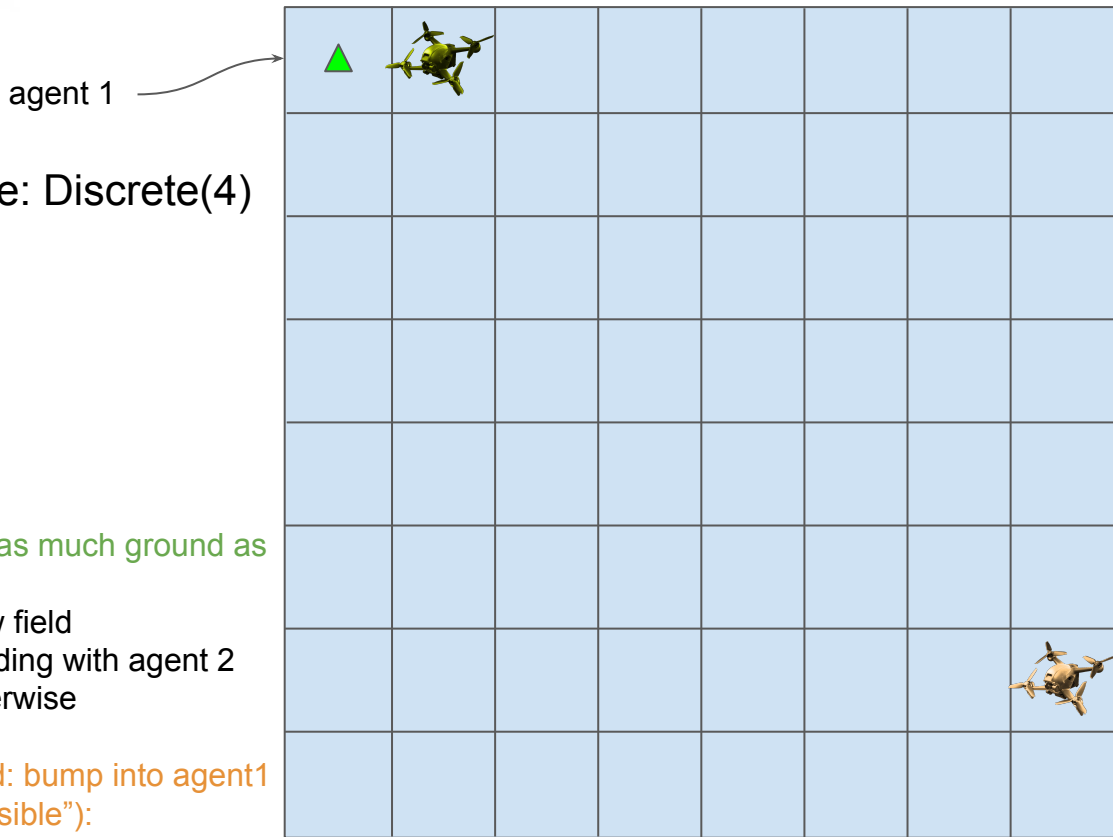MultiDiscrete([64, 64])

pos. agent 2

pos. agent 1

agent 2

# Our Environment - The "MultiAgentArena" 🫥



action space: Discrete(4)
0 - up
1 - right
2 - down
3 - left

rewards:
agent1 ("cover as much ground as possible"):
- +1 if new field
- -1 if colliding with agent 2
- -0.5 otherwise

agent2 ("defend: bump into agent1 as often as possible"):
- +1 if colliding with agent 1
- -0.1 otherwise

agent 1

observation space:
MultiDiscrete([64, 64])

pos. agent 2

pos. agent 1

agent 2

# Our Environment - The "MultiAgentArena" 🙂



agent 1

action space: Discrete(4)
0 - up
1 - right
2 - down
3 - left

rewards:
agent1 ("cover as much ground as possible"):
- +1 if new field
- -1 if colliding with agent 2
- -0.5 otherwise

agent2 ("defend: bump into agent1 as often as possible"):
- +1 if colliding with agent 1
- -0.1 otherwise

observation space:
MultiDiscrete([64, 64])

pos. agent 2

pos. agent 1

agent 2

# Our Environment - The "MultiAgentArena" 🫥



agent 1

action space: Discrete(4)
0 - up
1 - right
2 - down
3 - left

observation space:
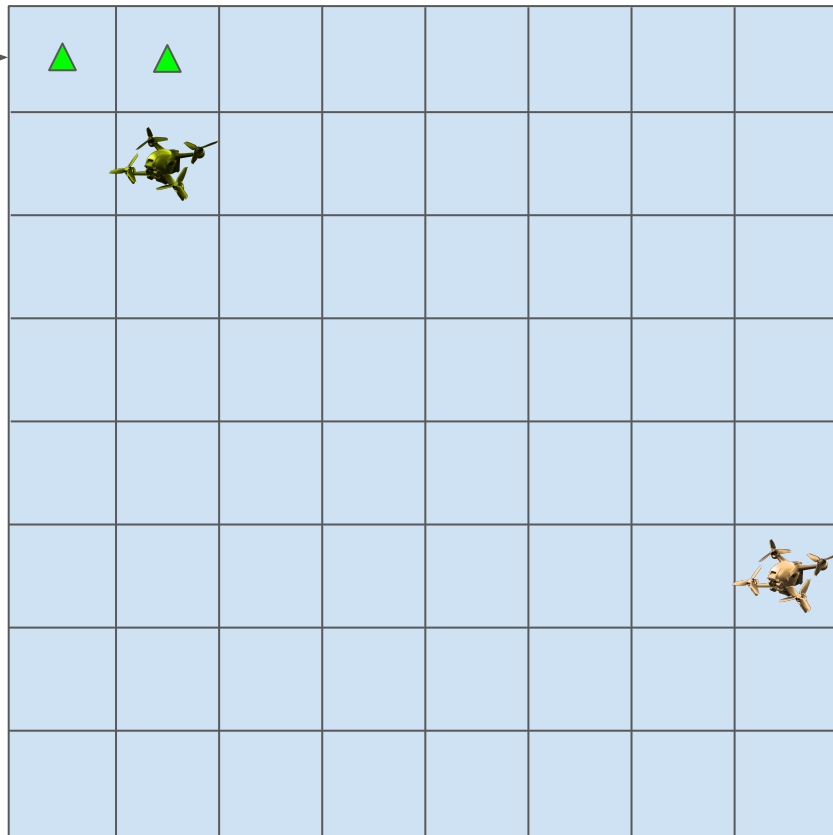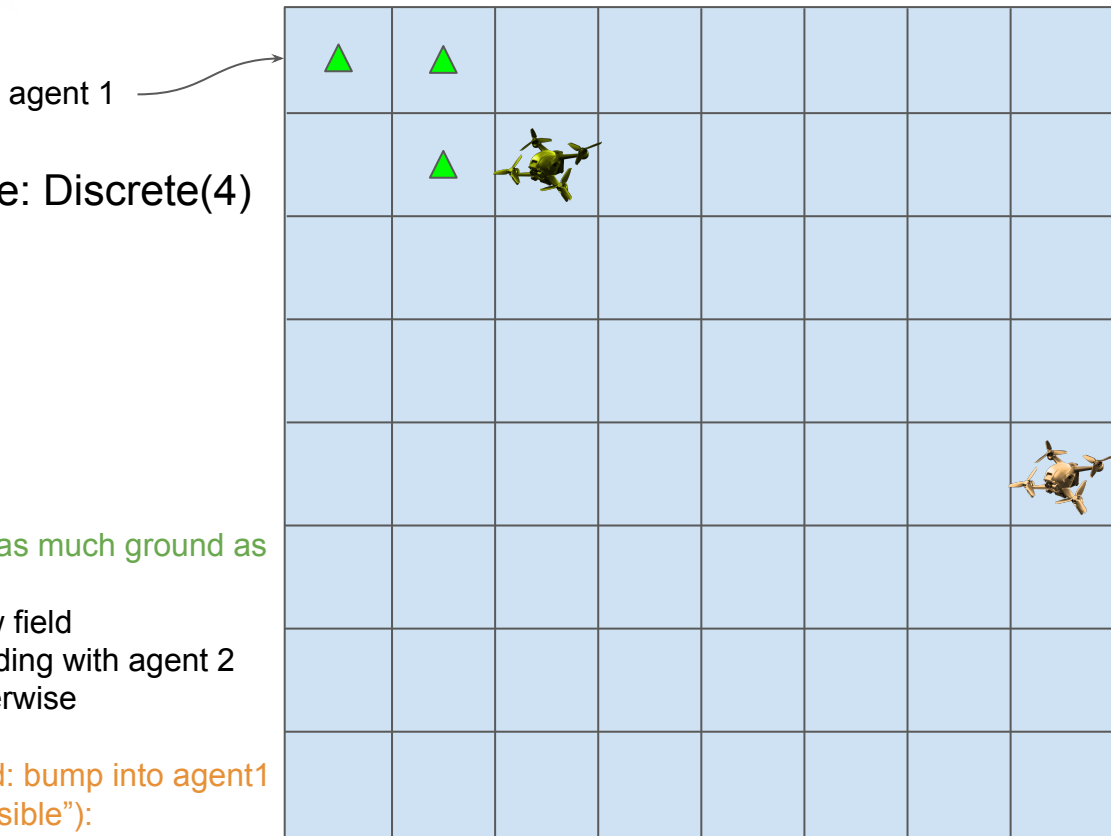MultiDiscrete([64, 64])

pos. agent 2

pos. agent 1

rewards:
agent1 ("cover as much ground as possible"):
- +1 if new field
- -1 if colliding with agent 2
- -0.5 otherwise

agent2 ("defend: bump into agent1 as often as possible"):
- +1 if colliding with agent 1
- -0.1 otherwise

agent 2

# Our Environment - The "MultiAgentArena" 🫥

agent 1

action space: Discrete(4)
0 - up
1 - right
2 - down
3 - left

observation space:
MultiDiscrete([64, 64])

pos. agent 2

pos. agent 1

rewards:
agent1 ("cover as much ground as possible"):
- +1 if new field
- -1 if colliding with agent 2
- -0.5 otherwise

agent2 ("defend: bump into agent1 as often as possible"):
- +1 if colliding with agent 1
- -0.1 otherwise

agent 2

# Our Environment - The "MultiAgentArena" 😶



action space: Discrete(4)
0 - up
1 - right
2 - down
3 - left

observation space:
MultiDiscrete([64, 64])
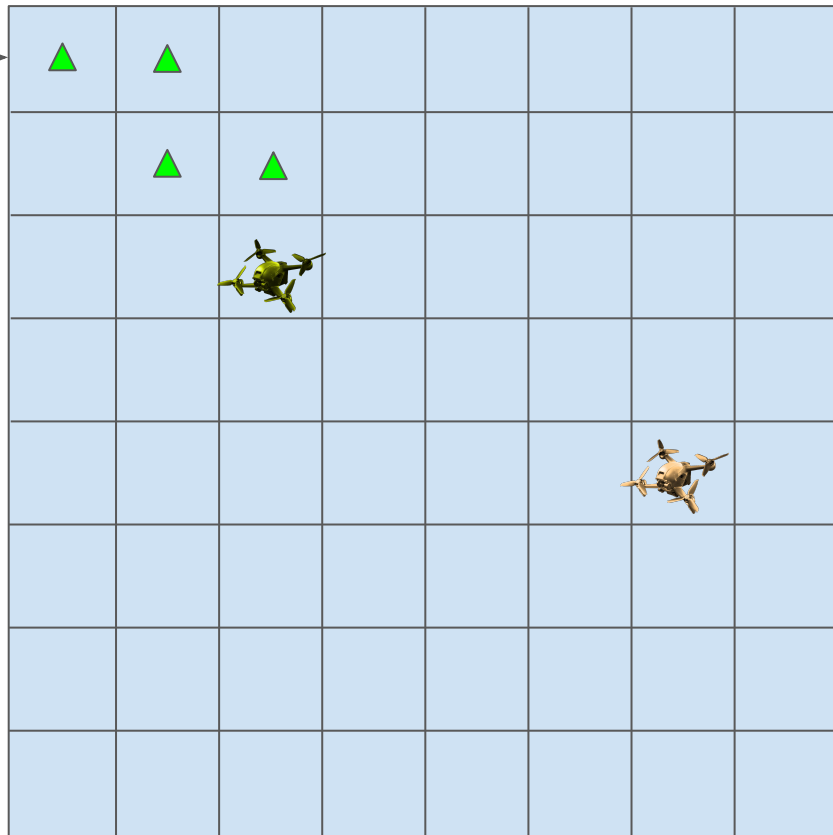
pos. agent 2

pos. agent 1

agent 1

agent 2

rewards:
agent1 ("cover as much ground as possible"):
- +1 if new field
- -1 if colliding with agent 2
- -0.5 otherwise

agent2 ("defend: bump into agent1 as often as possible"):
- +1 if colliding with agent 1
- -0.1 otherwise

# Our Environment - The "MultiAgentArena" 🫥

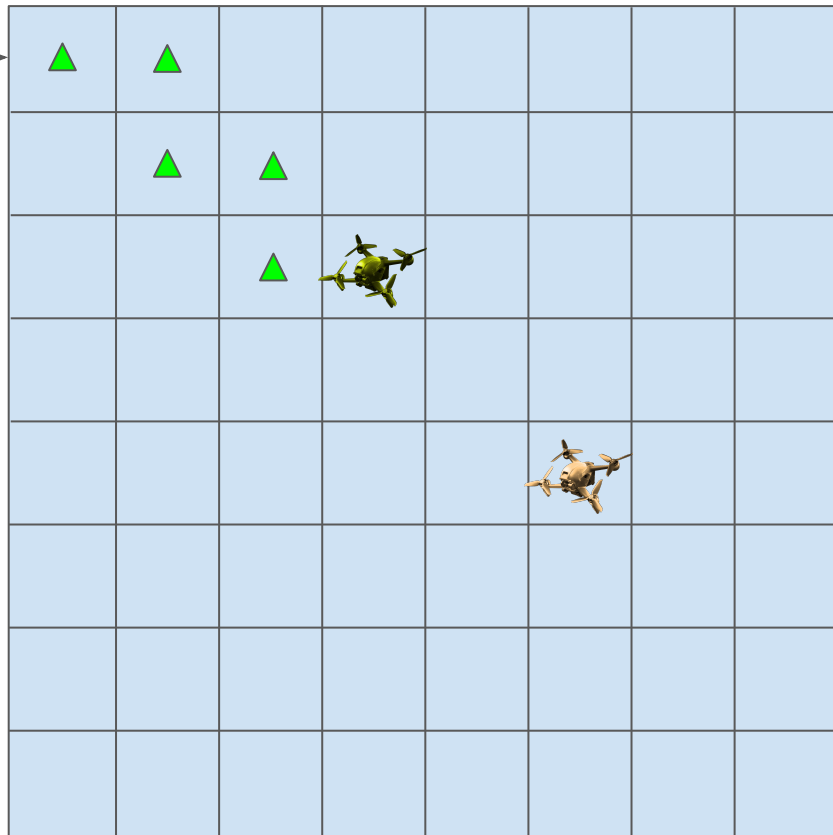action space: Discrete(4)
0 - up
1 - right
2 - down
3 - left

rewards:
agent1 ("cover as much ground as possible"):
- +1 if new field
- -1 if colliding with agent 2
- -0.5 otherwise

agent2 ("defend: bump into agent1 as often as possible"):
- +1 if colliding with agent 1
- -0.1 otherwise

agent 1

agent 2

observation space:
MultiDiscrete([64, 64])

pos. agent 2

pos. agent 1

# Our Environment - The "MultiAgentArena" 🙂

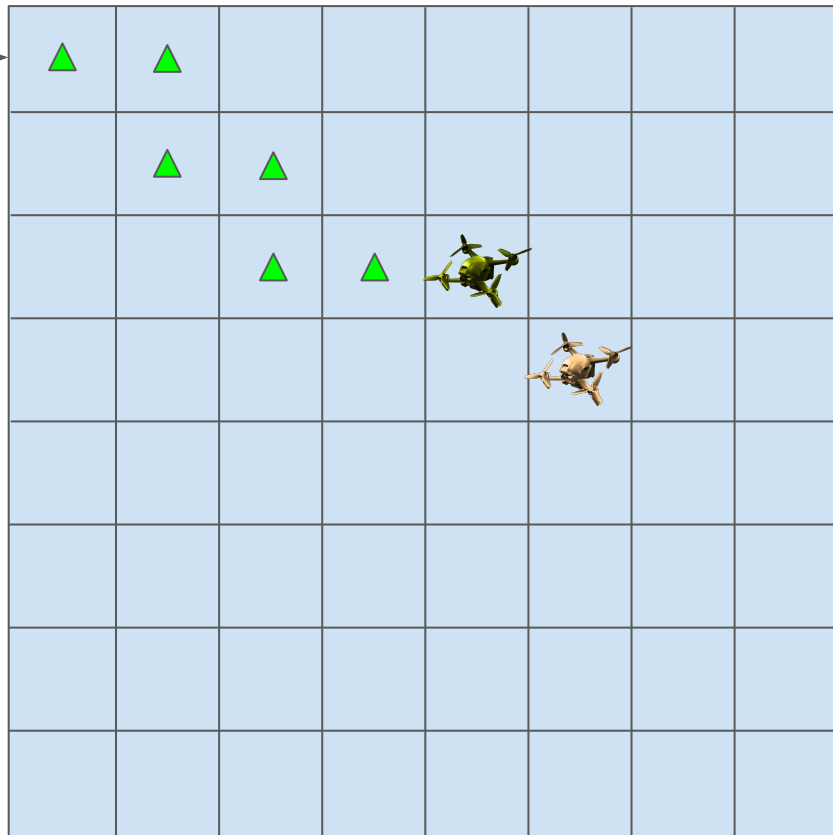action space: Discrete(4)
0 - up
1 - right
2 - down
3 - left

rewards:
agent1 ("cover as much ground as possible"):
- +1 if new field
- -1 if colliding with agent 2
- -0.5 otherwise

agent2 ("defend: bump into agent1 as often as possible"):
- +1 if colliding with agent 1
- -0.1 otherwise
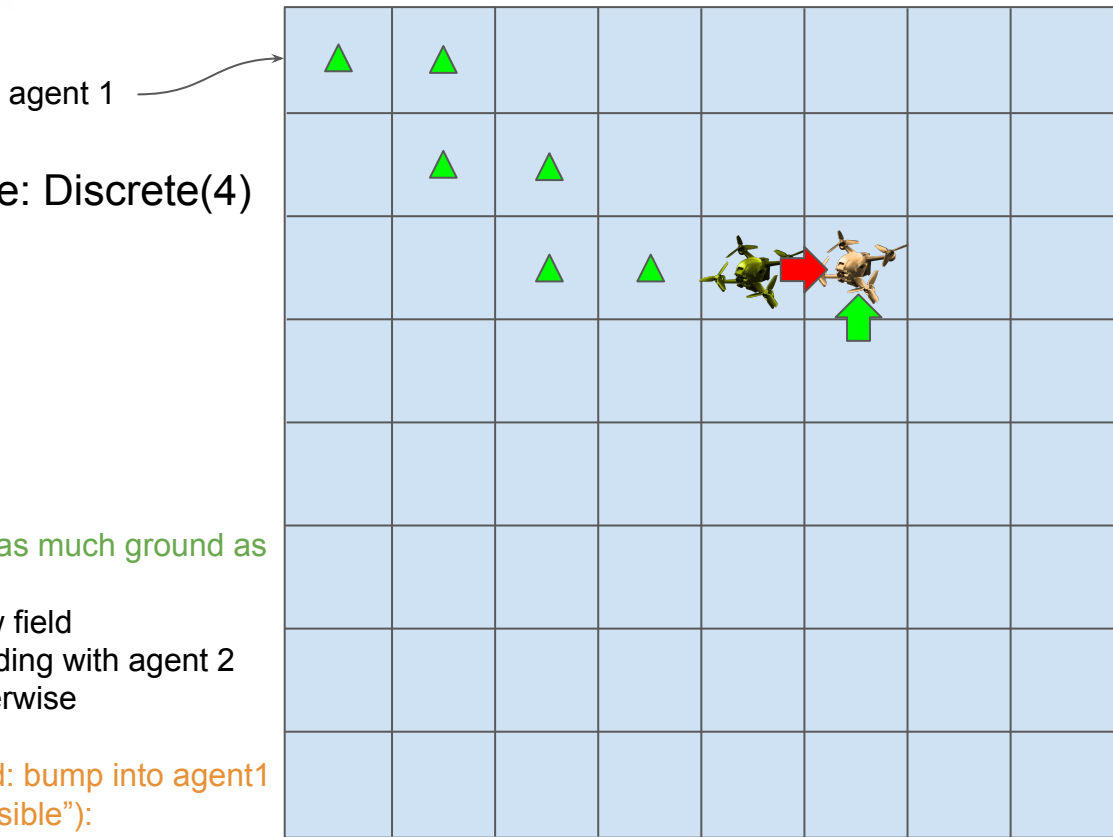
agent 1

observation space:
MultiDiscrete([64, 64])

pos. agent 2

pos. agent 1

agent 2

# What's a Space?



discrete
n=3

our action space

multi discrete
nvec=(3, 2, 8)

our obs. space

vector
shape=(4,)

matrix
shape=(4, 3)

4 dimensions
3 dimensions

tensor
shape=(4, 3, 3)

4 dimensions
3 dimensions
3 dimensions

Tuple space
shape=((4,), (3, 3, 2), (4, 3, 2))

Dict space
shape=((3, 3, 2), (3,))

"key Z":

"key A":

# Algo Overview

| Algorithm | Frameworks | Discrete Actions | Continuous Actions | Multi-Agent | Model Support | Multi-GPU |
|---|---|---|---|---|---|---|
| DQN, Rainbow | tf + torch | Yes +parametric | No | Yes | | tf + torch |
| APEX-DQN | tf + torch | Yes +parametric | No | Yes | | torch |
| IMPALA | tf + torch | Yes +parametric | Yes | Yes | +RNN, +LSTM auto-wrapping, +Attention, +autoreg | tf + torch |
| MAML | tf + torch | No | Yes | No | | torch |
| MARWIL | tf + torch | Yes +parametric | Yes | Yes | +RNN | torch |
| MBMPO | torch | No | Yes | No | | torch |
| PG | tf + torch | Yes +parametric | Yes | Yes | +RNN, +LSTM auto-wrapping, +Attention, +autoreg | tf + torch |
| PPO, APPO | tf + torch | Yes +parametric | Yes | Yes | +RNN, +LSTM auto-wrapping, +Attention, +autoreg | tf + torch |
| R2D2 | tf + torch | Yes +parametric | No | Yes | +RNN, +LSTM auto-wrapping, +autoreg | torch |
| SAC | tf + torch | Yes | Yes | Yes | | torch |

| Algorithm | Frameworks | Discrete Actions | Continuous Actions | Multi-Agent | Model Support | Multi-GPU |
|---|---|---|---|---|---|---|
| A2C, A3C | tf + torch | Yes +parametric | Yes | Yes | +RNN, +LSTM auto-wrapping, +Attention, +autoreg | A2C: tf + torch |
| ARS | tf + torch | Yes | Yes | No | | No |
| BC | tf + torch | Yes +parametric | Yes | Yes | +RNN | torch |
| CQL | torch | No | Yes | No | +RNN, +LSTM auto-wrapping, +autoreg | torch |
| ES | tf + torch | Yes | Yes | No | | No |

| | Continuous Actions | Multi-Agent | Model Support |
|---|---|---|---|
| | No | Yes | +RNN |
| | Partial | Yes | |

| | | |
|---|---|---|
| Shared Critic Methods | Depends on bootstrapped algorithm | |

Exploration-based plug-ins (can be combined with any algo)

| Algorithm | Frameworks | Discrete Actions | Continuous Actions | Multi-Agent | Model Support |
|---|---|---|---|---|---|
| Curiosity | tf + torch | Yes +parametric | No | Yes | +RNN |

A Closer Look at RLlib