



**FACULTY  
OF MATHEMATICS  
AND PHYSICS**  
Charles University

**MASTER THESIS**

Bc. Přemysl Bašta

# **Cooperation with Unknown Agents in Multi-agent Environment**

Department of Theoretical Computer Science and Mathematical Logic

Supervisor of the master thesis: : Mgr. Martin Pilát, Ph.D.

Study programme: Computer Science

Study branch: Artificial Intelligence

Prague 2023

I declare that I carried out this master thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In ..... date .....  
Author's signature

First of all, I would like to thank my beloved wife, Adéla, who carried me through most of my master’s studies and especially during the writing of this thesis, provided me with endless words of motivation and encouragement, and created an overall nurturing work environment. I would also like to thank my mother for her lifelong encouragement.

A big thank you goes to my wonderful supervisor Martin, who not only actively supported me during the time we worked on this thesis, but also inspired me to pursue nature-inspired science and multi-agent systems in the first place. I would also like to thank Milan Straka, whose amazing courses on Deep Learning and Deep Reinforcement Learning made it possible for me to dive deeper into the field through this thesis.

Finally, I would like to thank the MetaCentrum (MetaVO) organization for providing me with cluster computing resources, which made all the necessary computations possible.

Title: Cooperation with Unknown Agents in Multi-agent Environment

Author: Bc. Přemysl Bašta

Department: Department of Theoretical Computer Science and Mathematical Logic

Supervisor: : Mgr. Martin Pilát, Ph.D., Department of Theoretical Computer Science and Mathematical Logic

Abstract: Over the past few decades, we have witnessed great successes in the field of deep and reinforcement learning. Great achievements have been made in many competitive settings, both in single-agent and multi-agent environments, where AI has managed to outperform human experts and even entire teams of human experts. However, the situation is much more difficult when cooperation is required in purely cooperative environments. We first give a brief overview of reinforcement learning theory and its current state of the art algorithms. We then extend the theory to multi-agent systems, where several related issues are discussed. And finally, we propose novel approaches of agent training where we use a simplified multi-agent cooperative cooking game environment based on the popular game Overcooked, we attempt to train agents that are robust and capable of ad hoc agent cooperation.

Keywords: Reinforcement Learning Multi-Agent Systems Ad-hoc Cooperation

# Contents

<b>Introduction</b>	<b>3</b>
<b>1 Introduction to reinforcement learning</b>	<b>5</b>
1.1 Markov decision process . . . . .	5
1.2 Single-agent environment . . . . .	5
<b>2 Reinforcement learning algorithms</b>	<b>10</b>
2.1 Q-Learning . . . . .	10
2.2 Policy gradient methods . . . . .	13
2.2.1 Idea . . . . .	13
2.2.2 Policy Gradient Theorem . . . . .	13
2.2.3 Vanilla Policy Gradient . . . . .	14
2.2.4 Trust Region Policy Optimization . . . . .	15
2.2.5 Proximal Policy Optimization . . . . .	17
<b>3 Multi-agent environment</b>	<b>21</b>
3.1 Multi-agent Markov Decision Process . . . . .	21
3.2 Decentralized Partially Observable MDP . . . . .	21
3.2.1 Nash Equilibria . . . . .	23
3.3 Learning schemes . . . . .	23
3.4 Non-stationarity . . . . .	24
3.5 RL algorithms . . . . .	25
<b>4 Overcooked environment</b>	<b>27</b>
4.1 Overcooked game . . . . .	27
4.2 Basic layouts . . . . .	27
4.3 Environment description . . . . .	29
4.3.1 Action space . . . . .	29
4.3.2 State representation . . . . .	29
4.3.3 Rewards . . . . .	30
4.3.4 Initial state . . . . .	30
<b>5 Related work</b>	<b>32</b>
5.1 Human cooperation . . . . .	32
5.2 AI-AI cooperation . . . . .	32
5.2.1 Self-play . . . . .	33
5.2.2 Partner Sampling . . . . .	33
5.2.3 Population-Based Training . . . . .	33
5.2.4 Pre-trained Partners . . . . .	33
5.2.5 Result . . . . .	33
5.3 Problem of robustness . . . . .	34
5.3.1 Average performance . . . . .	34
5.3.2 Threshold performance . . . . .	34
5.3.3 Edge case testing . . . . .	34

<b>6</b>	<b>Our work - Preparation</b>	<b>36</b>
6.1	Framework . . . . .	36
6.1.1	Modifications . . . . .	36
6.2	Self-play . . . . .	38
6.2.1	Random state initialization . . . . .	41
6.2.2	Focus on Forced coordination . . . . .	41
<b>7</b>	<b>Our work - Contribution</b>	<b>43</b>
7.1	Multi agent setting simplification . . . . .	43
7.1.1	Single-agent perspective . . . . .	43
7.1.2	Non-stationarity . . . . .	43
7.1.3	Population exploitance . . . . .	44
7.2	Population building . . . . .	44
7.2.1	Inicialization . . . . .	44
7.2.2	KL divergence . . . . .	44
7.2.3	Population structure . . . . .	47
7.3	Simple convolution experiments . . . . .	47
7.3.1	Populations evaluation . . . . .	51
	<b>Conclusion</b>	<b>54</b>
	<b>Bibliography</b>	<b>55</b>

# Introduction

In recent years, people have witnessed rapid progress in artificial intelligence (AI) in all kinds of fields. Beating world champions at chess or Go is no longer a problem for AI models. The same pattern can be seen in more recent popular games such as Dota (OpenAI [2019]) or Starcraft, where even great human-AI team cooperation behavior has been achieved. However, all of these examples share the common trait of being competitive. The ultimate goal of our society is to create AI that cooperates with humans, not competes with them.

Recent work has shown that cooperative AI models trained together on purely cooperative tasks tend to rely on near-optimal behavior from their partners, and fail to cooperate with partners who don't meet this condition. This is bad news for us humans, because our behavior is rarely optimal.

A great example of human-AI cooperation where humans do not always perform perfectly are self-driving cars. In a situation where an accident is imminent, humans have to react quickly without having enough time to consider all possible reactions or even analyze the entire current road situation. However, car accidents are perhaps even too extreme an example of human suboptimal behavior. Nevertheless, people often fail at the even simpler task of following standard traffic rules when they have enough time to react. We can imagine that predicting human behavior is not an easy task for a self-driving car.

In this work, we will first operate in a single-agent environment, revisiting the definition of Markov decision, the building block of reinforcement learning. Based on this theory, we will intuitively introduce popular reinforcement learning algorithms divided into two categories of Q-learning and policy optimization. We primarily focus on policy branch of reinforcement algorithms, especially policy learning algorithm proximal policy optimization, which is considered as state of the art algorithm massively deployed in many successful projects.

Subsequently, we extend the theory developed for single-agent environments to more complex scenarios where more agents are involved. We highlight problems related to multi-agent settings, where the observability of the world is often an issue compared to single-agent environments. We illustrate problems with multi-agent training, where changes in one agent affect the behavior of the environment from the point of view of other agents, introducing the problem of non-stationarity. Finally, some single-agent variants of reinforcement learning algorithms are extended to multi-agent settings, where we mention the important aspects of these extensions.

We will use a simplified cooperative cooking environment based on the popular video game Overcooked, where two partners are forced to coordinate a shared task of cooking and delivering soup to a customer. We will familiarize ourselves with several different kitchen layouts, as each layout may offer different cooperative obstacles. Here we summarize what approaches have been tested in related work with respect to ad hoc agent cooperation. We mention the problem of defining the robustness of agent cooperation and different possible definitions of robustness.

And in the last part of this work, we prepare our working environment by modifying the stable-baselines3 library designed for reinforcement learning algorithms. We try to reimplement some methods of previously related work regarding

the problem of ad hoc coordination in AI-AI settings, both for verification purposes and also for building us an evaluation tool for our experiments. And finally, we propose a diversification method for building a population of agents that are designed to try to differentiate their behavior from those they have encountered in the training population. After tuning and optimizing our experiments on a selected kitchen layout, we then evaluate our approach on some of the remaining layouts.



# 1. Introduction to reinforcement learning

In this introductory chapter, we slowly build up the intuition and motivation behind reinforcement learning. We start by defining the mathematical model of the Markov decision process and then proceed with other related properties and relationships. The following pages are inspired by the introduction to reinforcement learning as presented by the authors of the spinningup library (Achiam [2018]) and also by the introductory book to reinforcement learning (Sutton and Barto [2018]).

## 1.1 Markov decision process

### Definition

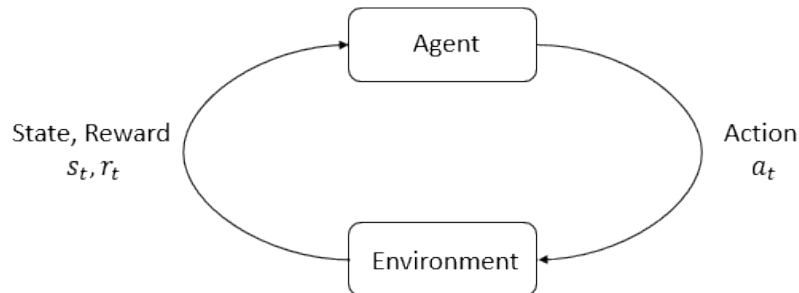
**Markov Decision Process** is tuple  $\langle S, A, R, P, \rho_0 \rangle$ , where

- $S$  is the set of all valid states,
- $A$  is the set of all valid actions,
- $R : S \times A \times S \rightarrow \mathbb{R}$  is the reward function, with  $r_t = R(s_t, a_t, s_{t+1})$  being reward obtained when transitioning from state  $s_t$  to  $s_{t+1}$  using action  $a_t$ .
- $P : S \times A \rightarrow \mathcal{P}(S)$  is the transition probability function, where  $P(s_{t+1}|s_t, a_t)$  is probability of transition from state  $s_t$  to  $s_{t+1}$  after taking action  $a_t$ .
- $\rho_0$  is starting state distribution.

The name Markov comes from the fact that the system satisfies the Markov property, which states that the history of previous states has no effect on the next state and that only the current state is considered for state transitions.

## 1.2 Single-agent environment

Having defined the mathematical model of the environment, let's review the related concepts. The whole problem of reinforcement learning (RL) can be best described by the following visualization.



Environment represents a kind of world with its internal rules and properties. The agent is then an entity that exists within this world, observes **state**  $\mathbf{s}$  of the world, decides to react with **action**  $\mathbf{a}$  on the basis of this state, and receives **reward**  $\mathbf{r}$  as a consequence of this action. The entire mechanism of this environment can then be broken down into these cycles of states, actions, and rewards. The goal of an agent is to interact with the environment in such a way as to maximize its cumulative reward.

## Observability

State  $\mathbf{s}$  contains all information about environment at given time. However, in some environment agent can perceive only **observation**  $\mathbf{o}$  where some information about environment can be missing. In this case we say that environment is **partially observable** as opposed to **fully observable** environment where agent has all information available at it's observation. Nevertheless, this problem is more related to multi-agent environments discussed in the next chapter and not so much to single-agent settings.

## Actions and policies

Environments can also differ in terms of what actions are possible within a given world. The set of possible actions is called **action space**, which again can be divided into two types. **Discrete** action space contains finite number of possible actions. And **continuous** action space, which allows the action to be any real-valued number or vector.

The agent's choice of action can then be described by a rule called **policy**. A common notation is that if the action selection is deterministic, we say the policy is **deterministic** and denote by

$$a_t = \mu(s_t).$$

If policy is **stochastic** it is usually noted as

$$a_t \sim \pi(\cdot | s_t).$$

Policies are the main object of interest of reinforcement learning, as this action selection mechanism of an agent is what we are trying to learn. A policy, for optimization purposes, is a function often parametrized by a neural network whose parameters are usually denoted by the symbol  $\theta$ , therefore, parameterized deterministic and stochastic policies are represented by the symbols  $\mu_\theta(s_t)$ ,  $\pi_\theta(\cdot | s_t)$  respectively.

## Trajectory

The next important definition is the notion of trajectory, also known as episode or rollout. A trajectory is a sequence of states and actions in an environment.

$$\tau = (s_0, a_0, s_1, a_1, \dots)$$

The initial state  $s_0$  of an environment is sampled from **start-state distribution**, denoted as  $\rho_0$ . Subsequent states follow the transition function of the environment. These may again be deterministic

$$s_{t+1} = f(s_t, a_t)$$

or stochastic,

$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

## Return

We have already mentioned the agent's desire to maximize cumulative rewards. Now we combine it with trajectories and derive the formulation of **return**.

$$R(\tau) = \sum_{t=0}^{|\tau|} r_t \quad (\text{finite-horizon})$$

$$R(\tau) = \sum_{t=0}^{|\tau|} \gamma^t r_t \quad (\text{infinite-horizon discounted return})$$

Infinite-horizon discounting is both intuitive and mathematically convenient.

## Optimal policy

In general, the goal of RL is to find such a policy that maximizes the expected return when acted upon. Suppose both the environment state transitions and the policy are stochastic. Then we can define the probability of the trajectory as

$$P(\tau|\pi) = \rho_0(s_0) \prod_{t=0}^{|\tau|} P(s_{t+1}|s_t, a_t) \pi(a_t|s_t).$$

The expected return  $J(\pi)$  can then be expressed as

$$J(\pi) = \int_{\tau} P(\tau|\pi) R(\tau) = \mathbb{E}_{\tau \sim \pi} [R(\tau)].$$

And finally, we can finish by defining the optimal policy

$$\pi^* = \arg \max_{\pi} J(\pi)$$

which is also an expression describing the central RL optimization problem.

## Value functions

Once we have some policy  $\pi$  it would be useful to define value of observed state. For that matter we define two functions.

**On-Policy Value Function**  $V^{\pi}(s)$ , which yields value of expected return when starting from state  $s$  and following policy  $\pi$ :

$$V^{\pi}(s) = \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s]$$

Similarly we define **On-Policy Action-Value Function**  $Q^\pi(s, a)$  which adds the possibility to say that in state  $s$  we take an arbitrary action  $a$  that does not necessarily have to come from policy  $\pi$ :

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s, a_0 = a]$$

For the optimal policy we further define **optimal value function**  $V_{\pi^*(s)}$  and **optimal action-value function**  $Q_{\pi^*(s,a)}$ :

$$\begin{aligned} V^*(s) &= \max_{\pi} \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s], \\ Q^*(s, a) &= \max_{\pi} \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s, a_0 = a] \end{aligned}$$

## Bellman equations

There exist formulations called Bellman equations that provide a way how to express value function in terms of action-value function and vice versa. They are based on the idea that the value of a state is equal to the reward you get in a given state, plus the value of the state you will obtain in the next transition. This idea also provides recursive relation.

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_{a \sim \pi(s)} [Q^\pi(s, a)] \\ &= \mathbb{E}_{a \sim \pi(s), s' \sim P(\cdot | s, a)} [R(s, a, s') + \gamma V^\pi(s')] \\ Q^\pi(s, a) &= \mathbb{E}_{s' \sim P(\cdot | s, a)} [R(s, a, s') + \gamma V^\pi(s')] \\ &= \mathbb{E}_{s' \sim P(\cdot | s, a)} [R(s, a, s') + \gamma \mathbb{E}_{a' \sim \pi(s')} [Q^\pi(s', a')]] \end{aligned}$$

The most important theorem for us is the reformulation of Bellman's equations for optimal policies:

$$\begin{aligned} V^*(s) &= \max_a \mathbb{E}_{s' \sim P} [R(s, a, s') + \gamma V^*(s')] \\ Q^*(s, a) &= \mathbb{E}_{s' \sim P} [R(s, a, s') + \gamma \max_{a'} Q^*(s', a')] \end{aligned}$$

As we will see in the next section. The first RL algorithm will be a straightforward application of the Bellman equation for optimal policy.

## Advantage function

Now that we've spent a few sections defining functions for the absolute value of actions or state-action pairs, it's worth considering relative value as well. Often, when dealing with RL problems, it is not so important for us to know the exact value of the action-state pair, but rather whether and by how much a given action is, on average, relatively better than others. In other words, we want to know the relative advantage of a given action over others. For a given policy  $\pi$ , the advantage function  $A^\pi(s, a)$  describes how much better it is to take action  $a$  over

randomly sampled actions following policy  $\pi$ , assuming that policy  $\pi$  is followed in all subsequent steps. Mathematically, the advantage function is defined as follows

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s).$$

The concept of an advantage function will be an integral part of policy gradient based methods, as we will see in a later section.

## 2. Reinforcement learning algorithms

In the previous chapter, we built the theoretical foundation of reinforcement learning. In this chapter, we translate that theory into practical algorithms. There are no precise classification boundaries between RL algorithms, as many of the techniques described are shared and crossed by all kinds of algorithms. Therefore, it is difficult to come up with an absolutely definitive taxonomy. However, for our purposes and the scope required, the following division into Q-learning and policy optimization is sufficient.

### 2.1 Q-Learning

#### Idea

We will start with a family of algorithms that focuses on learning the action value approximator  $Q_\theta(s, a)$ , as described in the previous chapter. For this reason, the group of such algorithms can also be referred to as Q-learning. Our primary goal in the RL problem is to find a policy that the agent can follow. In the case of Q-learning, once we have learned the approximator  $Q_\theta(s, a)$ , we can derive the policy by always taking the best possible action in the given state according to the learned action-value function.

$$a(s) = \arg \max_a Q_\theta(s, a).$$

By incorporating Bellman's optimal policy equations, we can directly train the Q-network by minimizing the loss

$$L(\theta) = (r + \gamma \max_{a'} Q_\theta(s', a') - Q_\theta(s, a))^2.$$

And by computing the loss gradient, we arrive at the update rule

$$Q_\theta(s, a) = Q_\theta(s, a) + \alpha(r + \gamma \max_{a'} Q_\theta(s', a') - Q_\theta(s, a)),$$

which is the backbone of the 2.1 algorithm bearing the same name.

#### Instability

The algorithm has rarely been used in this pure form. It has primarily been described for tabular methods, where the action-value function is represented by a table instead of a network approximator. In its simplest form, training is unstable and suffers from a number of significant shortcomings. Most notable is the theoretical deadly triad counter example Sutton and Barto [2018], which consists of a combination of value approximation, bootstrapping, and off-policy training that can lead to instability and divergence. The value approximation condition is met because we use a Q-network to approximate the action value. Bootstrapping means that the estimate is used to compute the targets, this is also true in Q-learning for the same reason. Finally, the term off-policy stands for an approach where training data is collected using a different distribution than that of a target policy.

---

**Algorithm 2.1:** Q-learning

---

**Input:** initial action-value approximator Q parameters  $\theta$  .

1 **repeat**

2     Observe state  $s$  and select action  $a$  according to  $\epsilon$ -greedy w.r.t.  $Q$  e.g.

$$a = \begin{cases} \text{random action,} & \text{with probability } \epsilon, \\ \arg \max_a Q(s, a), & \text{otherwise.} \end{cases}$$

3     Execute  $a$  in the environment.

4     Observe next state  $s'$ , reward  $r$  and done signal  $d$  to indicate whether  $s'$  is terminal.

5     **if**  $d$  is true **then**

6         Reset environment state.

7     **end if**

8     Compute targets

9

$$y(r, s', d) = r + \gamma(1 - d) \max_{a'} Q_{\theta}(s', a')$$

10     Update Q-network taking one step of gradient decent on

$$(y(r, s', d) - Q_{\theta}(s', a))^2$$

11 **until** *convergence*;

---

## DQN

One of the most outstanding papers based on Q-learning was the algorithm Deep Q-learning 2.2, which demonstrated super-human results on several Atari games Mnih et al. [2015]. We give the pseudocode of the algorithm in it's original form. The notation may seem a bit different from ours, but it represents the same mechanisms that we expect. To address the problem of correlated transition sequences of data, they introduce experience replay, where previously sampled transitions are stored. During training, data is sampled from this buffer, thus smoothing the training distribution over different past behaviors. However, probably the most important idea was the usage of a target Q-network, which broke the value approximation condition of the deadly triad, thus making the algorithm more robust. The target network is a copy of the original Q-network, with its parameters frozen and updated only once in a while based on the parameters of the main network. Its sole purpose is to compute target estimates that are not directly dependent on the Q-network function.

## Rainbow

Deep Q Learning was a significant contribution that led to the study of further Q Learning capabilities. The Rainbow project Hessel et al. [2017] could probably be called the pinnacle of such research. In this paper, they further investigate several isolated ideas of possible improvements and try to combine them. To name a few, they use double Q-network to address the problem of maximization

---

**Algorithm 2.2:** Deep Q-learning with experience replay

---

```
1 Initialize replay memory  $D$  to capacity  $N$ 
2 Initialize action-value function  $Q$  with random weights  $\theta$ 
3 Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$ 
4 for  $episode = 1, M$  do
5   Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$ 
6   for  $t=1, T$  do
7     With probability  $\epsilon$  select a random action  $a_t$ 
8     otherwise select  $a_t = \arg \max_a Q(\phi(s_t), a; \theta)$ 
9     Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
10    Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
11    Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$ 
12    Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$ 
13    Set
        
$$y_j = \begin{cases} r_j & \text{if terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$$

14    Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with
        respect to the network parameters  $\theta$ 
15    Every  $C$  steps reset  $\hat{Q} = Q$ 
16 end for
```

---



bias and enhance sampling from experience buffer by considering the priority of stored individual data samples. Together with all the other improvements, they achieved state-of-the-art performance on the Atari 2600 benchmark, both in terms of data efficiency and final performance.

## 2.2 Policy gradient methods

### 2.2.1 Idea

In the previous section, we were acquainted with the first group of RL algorithms, where we derived the final policy by acting according to *argmax* of our Q-function approximator. Since our objective is to find an optimal policy, this approach of considering Q values may seem a bit indirect. Fortunately, there is a whole other family of algorithms that deal with this very issue. As the name suggests, policy gradient algorithms focus on directly optimizing the policy  $\pi_\theta(a|s)$ . This is achieved by directly taking steps along the gradient of the expected return  $J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[R(\tau)]$ . The optimization step then has the form

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi_\theta)|_{\theta_k}$$

and the gradient  $\nabla_\theta J(\pi_\theta)$  is called **policy gradient**.

Before converting this into an algorithm, we have to figure out how to compute the policy gradient numerically. Such an expression can be obtained as a result of the Policy Gradient Theorem (PGT).

TODO: Nasledující veta a důkaz je inspirovány odtud [https://spinningup.openai.com/en/latest/spinningup/rl\\_intro3.html](https://spinningup.openai.com/en/latest/spinningup/rl_intro3.html) kde se odkazují na původní literaturu, kde je skutečně PGT dokazován. Ale je tam pouze v Sumární podobě s hodně předpoklady a celkově dost jinou reprezentací. Nikde se tam veta nebo důkaz pro integrál nevyskytuje. Dává mi intuitivně dost smysl důkaz přepsat z diskretního stavu do spojitěho a tvrdit se že vše platí jako předtím, ale nevím, nevím, nepodařilo se mi najít nikde jinde odkaz na výsledek v této podobě. Zajímavé bylo když jsem se koukal do diplomky od Honzy Uhlíka, že on to tam uvádí také v integrační podobě, zavádí si tam ještě distribuci pro počáteční stav a do znení věty rovnou zahrne i fakt že expected return lze nahradit za jinou náhodnou proměnou nezávislou na akci - také tam má pouze odkaz na tu samou původní literaturu. Tak nevím, jak se k tomu postavit. Jen tak prohlásit že když to platí pro diskretní podmínky, tak že to platí i pro spojitě, mi připadá zvláštní.

### 2.2.2 Policy Gradient Theorem

Policy Gradient Theorem (Sutton and Barto [2018]). It holds:

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{|\tau|} \nabla_\theta \log \pi_\theta(a_t|s_t) R(\tau) \right]$$

This can be proven by rewriting the formula in the following way:

$$\begin{aligned}
\nabla_{\theta} J(\pi_{\theta}) &= \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] \\
&= \nabla_{\theta} \int_{\tau} P(\tau|\theta) R(\tau) \\
&= \int_{\tau} \nabla_{\theta} P(\tau|\theta) R(\tau) \\
&= \int_{\tau} P(\tau|\theta) \nabla_{\theta} \log P(\tau|\theta) R(\tau) \\
&= \mathbb{E}_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log P(\tau|\theta) R(\tau)] \\
&= \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^{|\tau|} \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) R(\tau) \right]
\end{aligned}$$

It is useful to have a look at the expression and what it represents. There are two important components of the expression:  $\pi_{\theta}(a_t|s_t)$  and  $R(\tau)$ . Taking the gradient step of this objective, we are actually stating that we want to make the change in log probability  $\pi_{\theta}(a_t|s_t)$  weighted by how good the expected return was. However, this may feel slightly counterintuitive, since the expected return takes into account all the rewards of the episode. We may want to restrict ourselves to the future consequences of a given action. Fortunately, it can be shown that  $R(\tau)$  can be replaced by many other useful functions (Schulman et al. [2015b]):

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^{|\tau|} \Psi_t \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \right],$$

where  $\Psi_t$ , can be one of the following:

$\sum_{t=0}^{|\tau|} r_t$ : total reward of the trajectory

$\sum_{t=t'}^{|\tau|} r_t$ : reward following action  $a_{t'}$

$\sum_{t=t'}^{|\tau|} r_t - b(s_t)$ : baselined version of previous formula

$Q^{\pi}(s_t, a_t)$ : state-action value function

$A^{\pi}(s_t, a_t)$ : advantage function

### 2.2.3 Vanilla Policy Gradient

Common practice for policy gradient algorithms is to use some form of advantage function, where the baseline function is the value approximator  $b(s_t) = V^{\pi}(s_t)$ . The value approximator is typically represented by another neural network and is learned in parallel with the policy. There is a naming convention for the two types of networks. The policy network is usually called an actor since its job is to provide a policy for the agents to act on. The value network, on the other hand, is often referred to as a critic, because it produces, in a sense, a critique of the value of the state. Incorporating the value approximator and the idea of the utility function reduces the variance in the sample estimation and leads to more stable and faster learning.

With all of this said, we present the first algorithm of this class.

---

**Algorithm 2.3:** Vanilla Policy Gradient Algorithm

---

- 1 Input: initial policy parameters  $\theta_0$ , initial value function parameters  $\phi_0$
- 2 **for**  $k = 0, 1, 2, \dots$  **do**
- 3     Collect set of trajectories  $D_k = \{\tau_i\}$  by running policy  $\pi_k = \pi(\theta_k)$  in the environment.
- 4     Compute rewards-to-go  $\hat{R}$ .
- 5     Compute advantage estimates,  $\hat{A}_t$  (using any method of advantage estimation) based on the current value function  $V_{\phi_k}$ .
- 6     Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|\mathcal{D}_k|} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_k} \hat{A}_t.$$

- 7     Compute policy update, either using standard gradient ascent,

$$\theta_{l+1} = \theta_k + \alpha_k \hat{g}_k,$$

or via another gradient ascent algorithm like Adam. Fit value function by regression on mean-squared error:

8

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T (V_{\phi}(s_t) - \hat{R}_t)^2,$$

typically via some gradient descent algorithm.

9 **end for**

---

Policy gradient methods differ from Q-learning in several ways.

First, the absolute value of the objective function we are optimizing cannot be interpreted in terms of performance outcome. In fact, taking a step in the gradient descent does not even guarantee an improvement in the expected return in general. On a given set of samples, a value of even  $-\infty$  can be achieved. However, the expected return of a changed policy would most likely be abysmal.

And second, policy gradient algorithms are **on-policy**, meaning that only data samples collected using the most recent policy are used for training. This is in contrast to the off-policy approach of Q-learning, where training data samples collected throughout the training process are used for learning steps. For this reason, policy gradient algorithms, especially vanilla policy gradient algorithms, are often considered sample-inefficient compared to off-policy algorithms.

## 2.2.4 Trust Region Policy Optimization

Although the standard policy gradient step makes a small policy change within the parameter space, it turns out that it may have a significant change on the performance difference. Therefore, vanilla policy gradient algorithm must be careful not to take large steps, making it even more sample-inefficient.

The next algorithm in the policy gradient family attempts to address this problem. As its name suggests, Trust Region Policy Optimization (Schulman et al. [2015a]) tries to take steps within the trusted region in which it is con-

strained so as not to degrade performance. TRPO proposes theoretical update of parameterized policy  $\pi_\theta$  as

$$\begin{aligned}\theta_{k+1} &= \arg \max_{\theta} \mathcal{L}(\theta_k, \theta) \\ \text{s.t. } &\bar{D}_{KL}(\theta || \theta_k) \leq \delta\end{aligned}$$

where

$$\mathcal{L}(\theta_k, \theta) = \mathbb{E}_{s, a \sim \pi_{\theta_k}} \left[ \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a) \right]$$

is a surrogate advantage measuring the performance of the policy  $\pi_\theta$  relative to the old policy  $\pi_{\theta_k}$ .

And

$$\bar{D}_{KL}(\theta || \theta_k) = \mathbb{E}_{s \sim \pi_{\theta_k}} [D_{KL}(\pi_\theta(\cdot|s) || \pi_{\theta_k}(\cdot|s))]$$

is the average KL-divergence (Kullback [1959]) between policies evaluated on states visited by the old policy. However, working with theoretical updates of TRPO in this form is not an easy task. Therefore, approximations obtained by applying Taylor expansion around  $\theta_k$  are being used:

$$\begin{aligned}\mathcal{L}(\theta_k, \theta) &\approx g^T(\theta - \theta_k) \\ \bar{D}_{KL}(\theta || \theta_k) &\approx \frac{1}{2}(\theta - \theta_k)^T H(\theta - \theta_k)\end{aligned}$$

And the original problem can then be reformulated as an approximate optimization problem:

$$\begin{aligned}\theta_{k+1} &= \arg \max_{\theta} g^T(\theta - \theta_k) \\ \text{s.t. } &\frac{1}{2}(\theta - \theta_k)^T H(\theta - \theta_k) \leq \delta\end{aligned}$$

Such an approximate reformulation can be solved analytically using methods of Lagrangian duality (Rockafellar [1970]), yielding the solution:

$$\theta_{k+1} = \theta_k + \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g$$

However, since we employed Taylor expansion, the approximation error could violate the KL divergence constraint. For this reason, TRPO incorporates the idea of backtracking line search (Armijo [1966]):

$$\theta_{k+1} = \theta_k + \alpha^j \sqrt{\frac{2\delta}{g^T H^{-1} g}} H^{-1} g$$

where  $\alpha \in (0, 1)$  is the backtracking coefficient and  $j$  is the smallest non-negative integer such that the KL divergence constraint is fulfilled and the surrogate advantage is positive.

---

**Algorithm 2.4:** Trust Region Policy Optimization

---

- 1 Input: initial policy parameters  $\theta_0$ , initial value function parameters  $\phi_0$
- 2 Hyperparameters: KL-divergence limit  $\delta$ , backtracking coefficient  $\alpha$ , maximum number of backtracking steps  $K$
- 3 **for**  $k = 0, 1, 2, \dots$  **do**
- 4     Collect set of trajectories  $D_k = \{\tau_i\}$  by running policy  $\pi_k = \pi(\theta_k)$  in the environment.
- 5     Compute rewards-to-go  $\hat{R}$ .
- 6     Compute advantage estimates,  $\hat{A}_t$  (using any method of advantage estimation) based on the current value function  $V_{\phi_k}$ .
- 7     Estimate policy gradient as

$$\hat{g}_k = \frac{1}{|D_k|} \sum_{\tau \in D_k} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) |_{\theta_k} \hat{A}_t.$$

- 8     Use the conjugate gradient algorithm to compute

$$\hat{x}_k \approx \hat{H}_k^{-1} \hat{g}_k,$$

where  $\hat{H}_k^{-1}$  is the Hessian of the sample average KL-divergence.

- 9     Compute the policy by backtracking line search with

$$\theta_{l+1} = \theta_k + \alpha^j \sqrt{\frac{2\delta}{\hat{x}_k^T \hat{H}_k \hat{x}_k}} \hat{x}_k,$$

where  $j \in \{0, 1, 2, \dots, K\}$  is the smallest value which improves the sample loss and satisfies the sample KL-divergence constraint.

- 10     Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|D_k|T} \sum_{\tau \in D_k} \sum_{t=0}^T (V_{\phi}(s_t) - \hat{R}_t)^2,$$

typically via some gradient descent algorithm.

- 11 **end for**
- 

### 2.2.5 Proximal Policy Optimization

Finally, we will conclude this chapter with last algorithm from family of policy gradient algorithms, which is considered in many aspects to be the state of the art algorithm. As the authors of this next algorithm state. With the leading contenders Q-learning, "vanilla" policy gradient methods, and trust region policy gradient methods, there is still room for the development of a method that is

- scalable - large models and parallel implementations
- data efficient
- robust - successful on a variety of problems without hyperparameters tuning

Q-learning is poorly understood and fails on many trivial problems. Vanilla policy gradient methods suffer from poor data efficiency and robustness. And trust region policy optimization is relatively complicated and not well suited to architectures including noise or parameter sharing.

Proximal policy optimization (Schulman et al. [2017]) aims at data efficiency and reliability of TRPO performance while using only first-order optimization. The authors propose a novel objective with clipped probability ratios, which provides a pessimistic lower bound on the performance of the policy.

Let  $r_t(\theta)$  denote the probability ratio

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}.$$

Then the "surrogate" objective of TRPO can be expressed again in the form:

$$L^{CPI}(\theta) = \hat{\mathbb{E}}_t \left[ \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t [r_t(\theta) \hat{A}_t]$$

CPI refers to conservative policy iteration (Kakade and Langford [2002]), where this objective was originally proposed. Maximizing  $L^{CPI}$  without any constraint would, as discussed in the previous section, lead to an excessively large policy update. Thus, the authors propose a new modified objective, called the clipped surrogate objective, which penalizes policy adjustments that move  $r_t(\theta)$  far from 1.

$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t \left[ \min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right],$$

where the value  $\epsilon = 0.2$  is empirically suggested. The objective can be intuitively explained as following. The first term in min is the previous  $L^{CPI}$ . And the second term modifies the surrogate objective by a clipping probability ratio that prevents  $r_t$  from escaping the interval  $[1 - \epsilon, 1 + \epsilon]$ . Finally, taking the minimum of the clipped and unclipped objectives makes the final objective pessimistically bounded on the unclipped objective.

The training process of PPO then proposes to alternately sample data from the policy and then perform several epochs of optimization steps on the sampled data. Note here that for each first training epoch it holds that  $r_t(\theta) = 1$ , so the objective in the first epoch always equals  $L^{CPI}$ .

Alternatively, the authors propose a second version of PPO that includes the KL penalty as part of the objective, making it even more similar to the idea proposed in TRPO. However, we will not cover more details here, as the authors themselves prefer the version including clipping of the surrogate objective.

When using a neural network architecture that shares parameters between the policy and the value function, a combined objective function that includes both the policy surrogate and the value function error term must be applied.

$$L_t^{CLIP+VF}(\theta) = \hat{\mathbb{E}}_t \left[ L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) \right],$$

where  $c_1$  is the coefficient and  $L_t^{VF}$  is a squared error loss  $(V_\theta(s_t) - V_t^{targ})^2$ .

---

**Algorithm 2.5:** Proximal Policy Optimization

---

- 1 Input: initial policy parameters  $\theta_0$ , initial value function parameters  $\phi_0$
- 2 **for**  $k = 0, 1, 2, \dots$  **do**
- 3     Collect set of trajectories  $D_k = \{\tau_i\}$  by running policy  $\pi_k = \pi(\theta_k)$  in the environment.
- 4     Compute rewards-to-go  $\hat{R}_t$ .
- 5     Compute advantage estimates,  $\hat{A}_t$  (using any method of advantage estimation) based on the current value function  $V_{\phi_k}$ .
- 6     Update the policy by maximizing the PPO-Clip objective:

$$\theta_{k+1} = \arg \max_{\theta} \frac{1}{|D_k|T} \sum_{\tau \in D_k} \sum_{t=0}^T \min \left( \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} A^{\pi_{\theta_k}}(s_t, a_t), g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right),$$

typically via stochastic gradient ascent with Adam.

- 7     Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|D_k|T} \sum_{\tau \in D_k} \sum_{t=0}^T (V_{\phi}(s_t) - \hat{R}_t)^2,$$

typically via some gradient descent algorithm.

- 8 **end for**
- 

Lastly, one problem we haven't explicitly addressed is that of exploitation and exploration. In general, during training we always try to strike a balance between exploiting learned experience and exploring new possibilities. In the case of Q-learning, this is most often done artificially by using the  $\epsilon$ -greedy approach, where with probability  $1 - \epsilon$  we exploit our knowledge by taking  $\arg \max_a$  action according to the  $Q$  value. And with probability  $\epsilon$  we explore by taking a random action.

Sampling according to the policy  $\pi_{\theta}$  in policy optimization methods solves this problem more naturally. Typically, at the beginning of the training process, when the parameters are initialized, the parametrized probability distribution is close to uniform, which implicitly makes the action selection mechanism exploratory. And as the policy is updated, it becomes more specialized towards the optimal policy, effectively forcing the action selection to be more exploitative. Unfortunately, the exploration described in policy optimization methods is often insufficient, and policies tend to get stuck at bad local optima despite an initial uniform distribution. On this account, an exploration mechanism in the form of bonus entropy(Williams [1992]) is often employed and also recommended by the PPO authors. By adding a small bonus for policy entropy, the policy is slightly forced in the direction of uniform distribution, shifting towards exploratory behavior. Thus, the final PPO objective may have a form as follows:

$$L_t^{CLIP+VF}(\theta) = \hat{\mathbb{E}}_t \left[ L_t^{CLIP}(\theta) - c_1 L_t^{VF}(\theta) + c_2 S[\pi_{\theta}](s_t) \right],$$

where  $c_2$  is another coefficient and  $S$  denotes an entropy bonus.

## Q-Learning and Policy optimization

So far we have covered the simplest possible division rule between types of reinforcement learning algorithms. However, there are actually several algorithms that combine both Q-learning and policy optimization. In these algorithms, both the  $Q$  approximator and the policy approximator are learned simultaneously. To name just a handful of the more popular ones: Deep Deterministic Policy Gradient (DDPG, Lillicrap et al. [2015]), Twin Delayed DDPG (TD3, Fujimoto et al. [2018]), Soft Actor-Critic (SAC, Haarnoja et al. [2018]). We won't cover them in detail, as they are beyond the scope of our needs.



## 3. Multi-agent environment

In this chapter, we revisit concepts from the previous section and extend them to multi-agent settings. We introduce theoretical definitions for various types of settings and provide possible schemes for concurrent learning of multiple agents. And at the end, we mention two popular RL approaches in the field of MARL.

### 3.1 Multi-agent Markov Decision Process

So far, all the theory and algorithms that have been built have revolved around environments where there is a single agent. However, this is rarely the case in real-world problems. Much more often we encounter environments with multiple agents operating within them. In the general case, the agents are heterogeneous, which means that the agents may have different goals. Nevertheless, we will mostly focus on environments that are fully cooperative, meaning that the utility of any given state of the system is equivalent for all agents.

With this setting we can extend the definition from the MDP section 1.1:

#### Definition

**Multiagent Markov Decision Process**(Boutilier [1996]) is a tuple  $\langle n, S, \mathcal{A}, T, R \rangle$

- $n$  is the number of agents
- $S$  is the set of all valid states,
- $\mathcal{A}$  is the set of joint actions
- $T : S \times \mathcal{A} \times S \rightarrow [0, 1]$  is the transition function
- $R : S \rightarrow \mathbb{R}$  is a real-valued reward function, where reward determined by the reward function  $R(s)$  is received by the the entire collection of agents, or alternatively, all agents receive the same reward.

### 3.2 Decentralized Partially Observable MDP

Unfortunately, we cannot be satisfied with this definition, since MMDP provides a description that is far too idyllic for real-world problems. The MMDP model presumes that the state of the environment can be globally accessed by all agents. However, this is rarely the case in multi-agent environments. In general, the world can be of high complexity, and combined with limited sensory capabilities, the agent may perceive only limited observations describing only a part of the entire environment state. Therefore, we extend our model definition to include the concept of partial observability.

## Definition

**Decentralized partially observable Markov decision process** (Dec-POMDP, [Oliehoek et al. [2008]]) is tuple  $\langle n, S, \mathcal{A}, P, R, \mathcal{O}, O, h, b^0 \rangle$  where,

- $n$  is the number of agents
- $S$  is the set of all valid states,
- $\mathcal{A}$  is the set of joint actions
- $P : S \times \mathcal{A} \times S \rightarrow [0, 1]$  is the transition function
- $R : S \rightarrow \mathbb{R}$  is the immediate reward function
- $\mathcal{O}$  is the set of joint observations
- $O : \mathcal{A} \times S \rightarrow \mathcal{P}(\mathcal{O})$  is the observation function
- $h$  is the horizon of the problem
- $b^0 \in \mathcal{P}(S)$ , is the initial state distribution at time  $t = 0$

We define  $\mathcal{A} = \times_i \mathcal{A}^i$ , where  $\mathcal{A}^i$  is the set of actions available to agent  $i$ . Similarly,  $\mathcal{O} = \times_i \mathcal{O}^i$ , where  $\mathcal{O}^i$  is the set of observations available to agent  $i$ . Note that even this definitional extension does not provide us with a model capable of describing an environment where agents have different reward functions, which is needed for situations where agents have different goals or are even competing. This extension is possible with the definition of Stochastic game (Shapley [1953]). However, we don't need to provide the formal definition and extend our models, since the primary goal of this work will focus mainly on the cooperative setting where all agents have a common goal.

If the observation satisfies the condition that the individual observation of all agents uniquely identifies the true state of the environment, the environment is considered fully observable and such a Dec-POMDP can be reduced to MMDP.

## Notation

To denote common entities, we will use bold:  $\mathbf{a} = (a^1, \dots, a^n) \in \mathcal{A}$ . The common policy  $\boldsymbol{\pi}$  induced by the set of individual policies  $\{\pi^i\}_{i \in n}$  gives the mapping from states to common actions. Now we can use the similar notation as in the first chapter by using the bold symbols:

$$\text{Trajectory : } \boldsymbol{\tau} = (s_0, \mathbf{a}_0, s_1, \mathbf{a}_1, \dots)$$

$$\text{Return : } R(\boldsymbol{\tau}) = \sum_{t=0}^{\tau} r_t$$

$$\text{Probability of trajectory : } P(\boldsymbol{\tau}|\boldsymbol{\pi}) = \rho_0(s_0) \prod_{t=0}^{|\boldsymbol{\tau}|-1} P(s_{t+1}|s_t, \mathbf{a}_t) \boldsymbol{\pi}(\mathbf{a}_t|s_t)$$

$$\text{Expected return : } J(\boldsymbol{\pi}) = \int_{\boldsymbol{\tau}} P(\boldsymbol{\tau}|\boldsymbol{\pi}) R(\boldsymbol{\tau}) = \mathbb{E}_{\boldsymbol{\tau} \sim \boldsymbol{\pi}} [R(\boldsymbol{\tau})]$$

Given the joint utility function, it is useful to think of the collection of agents as a single agent whose goal is to produce an optimal joint policy. The problem with treating MMDP as a standard MDP where actions are distributed lies in coordination. In general, there are several different optimal joint policies. However, even if all agents choose their individual policies according to some optimal policy, there is no guarantee that they all pick from the same optimal joint policy. Such a final joint policy may be nowhere near the optimal one, and most likely even produce significantly worse performance. In theory, there are two simple ways to ensure optimal coordination. First, there could exist a central control mechanism 3.3 that can compute the joint policy and then communicate the result actions to all individual agents. Or second, each agent can communicate its choice of individual policy to the others. However, neither of these approaches are feasible in the real application.

### 3.2.1 Nash Equilibria

Alternatively, we can look at the MMDP from the perspective of an n-person game. Then the problem of determining an optimal joint policy can be viewed as a problem of optimal equilibrium selection. A Nash equilibrium (Nash [1950]) for a  $\pi^*$  can be defined as: A set of policies  $\pi^*$  is a Nash equilibrium if:

$$\forall i \in n, \forall \pi^i : J(\pi^{*-i}, \pi^i) \geq J(\pi^{*-i}, \pi^i)$$

However, not all Nash equilibria are optimal joint policies, as some may have lower utility than others. This makes multi-agent environments more sensitive to convergence to local suboptimal policies, as the only way to escape such an equilibrium is through coordinated modifications of all individual policies.

## 3.3 Learning schemes

### Centralized scheme

From a theoretical point of view, we could consider MMDP as an instance of single-agent MDP, where the goal would be to learn a central joint policy, which would then be distributed among individual agents. With this idea, we could solve MMDP problems using the same single-agent RL algorithms from the second chapter, only instead of learning a single policy, we would learn a joint policy. However, this approach has several shortcomings in real-world cases.

First, from a practical point of view, agents in this scenario would be somewhat passive entities whose job would only be to report perceived observations to some central authority. After all the observations have been collected by the central unit, a joint policy is constructed and distributed to the agents, who then blindly act as instructed by the central control. This has several serious problems. Imagine that some failure of the central mechanism occurs. Suddenly, the whole system of agents collapses because all agents lack individual autonomy.

Second, such an intensive communication between all agents and the central control mechanism may be too demanding or not even plausible for technical reasons.

Finally, from a theoretical point of view, the representation of such a joint policy grows exponentially with the number of agents with respect to both observations ( $\prod_{i=0}^n |\mathcal{O}^i|$ ) and actions ( $\prod_{i=0}^n |\mathcal{A}^i|$ ), which makes it unscalable.

## Concurrent scheme

At the other end of the spectrum is the concurrent scheme. Here, all kinds of global information are omitted and agents rely solely on their local observations. The training of such an agent is then completely independent.

## Centralized training with decentralized execution

And finally, somewhere in between the concurrent and centralized schemes, we can identify the so-called centralized training with decentralized execution. Here we consider the two life stages of agents. First, we train our agents in safe laboratory conditions, and only after the training is complete are they deployed in the real world.

During training, we can take advantage of the fact that we presumably have more information about the state of the environment at our disposal. Whether it is the true global state of the environment, local observations of other agents, or actions taken by others. All of this additional information can be incorporated into the training process to capture the true state to act upon. In other words, we want to make the training process as simple and accurate as possible.

Once the agents are deployed, they again depend solely on their local observations.

This learning scheme is often used in RL algorithms where both actors and critics are employed. For example, this is reflected in the PPO 2.5 algorithms mentioned above. During training, both actor and critic are trained in parallel, and the value estimate obtained by the critic can be used to provide more accurate information about the real value of the given state, making the actor's policy update more stable and accurate. Once training is complete, the agents, provided with local observations, are asked to take action based solely on the actor.

## 3.4 Non-stationarity

In addition to all the problems of partial observability, local equilibria, and policy distribution that have been mentioned so far there is another important problem that is not present in single-agent settings. Although we are able to reproduce the value function expressions:

$$\begin{aligned} V^\pi(s) &= \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s] \\ Q^\pi(s, \mathbf{a}) &= \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s, \mathbf{a}_0 = \mathbf{a}] \\ A^\pi(s, \mathbf{a}) &= Q^\pi(s, \mathbf{a}) - V^\pi(s) \end{aligned}$$

due to the non-stationarity problem, we cannot obtain the optimal value using Bellman equations as it was done in the single agent setting. Non-stationarity refers to the fact that changing policies of other agents as a consequence changes,

from a fixed agent’s perspective, the state transition and the reward function of the environment. Using the idea of Bellman equations for optimality is still possible. However, in a multi-agent setting, we lose the theoretical basis that promises convergence to the optimal policy, and training using the Q-learning update rule must be accompanied by some mechanisms to overcome the non-stationarity problem.

### 3.5 RL algorithms

Having extended our mathematical model to settings that satisfy the conditions for a multi-agent environment, we can continue with a brief mention of the multi-agent variants of reinforcement learning algorithms discussed in the previous chapter.

#### MADDPG

The first of the algorithms mentioned is the Multi Agent Deep Deterministic Policy Gradient (MADDPG, Lowe et al. [2017]) algorithm. The authors extend DDPG (briefly mentioned in 8) by using centralized learning with decentralized execution. To address the problem of non-stationarity, the authors propose sampling from the ensemble of policies for each agent to obtain more robust multi-agent policies. This algorithm has been extensively experimented in academia with respect to multi-agent coordination environments.

#### MAPPO

Finally, we want to mention the recent success in the form of Multi Agent Proximal Policy Optimization (MAPPO, Yu et al. [2021]) variant of the PPO 2.5 algorithm. The authors revisit the use of PPO in multi-agent settings. In recent years, MADDPG has usually been the first choice of algorithm when dealing with multi-agent environments, leaving PPO ommitted. The authors hypothesize that this is due to two reasons: first, the belief that PPO is less sample-efficient than off-policy methods, and second, the fact that common implementation and hyperparameter tuning techniques when using PPO in single-agent settings often do not yield strong performance when transferred to multi-agent settings.

The authors propose five important changes with respect to adaptation to multi-agent settings.

- Value normalization:  
employing value normalization using Generalized Advantage Estimation (Schulman et al. [2015b])
- Input Representation to Value Function:  
enhanced global state with agent-specific features are used for critic, where such representation does not have substantially higher dimension
- Training Data Usage:  
using smaller amount of trainig epochs to avoid non-stationarity problem and not splitting data to mini-batches to make policy update more accurate

- PPO Clipping:  
maintain a clipping ratio  $\epsilon$  under 0.2, within this range, tune  $\epsilon$  as trade-off between training stability and fast convergence
- PPO Batch Size:  
utilize a large batch size

Using these five concepts, they were able to achieve comparable or superior results to off-policy algorithms on several benchmark frameworks without any other domain or structural modifications.

## 4. Overcooked environment

### 4.1 Overcooked game

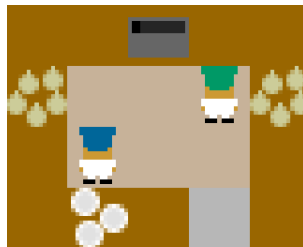
Before we get into our problems with cooperation, let us first examine the environment. We will work with an environment based on a popular cooking video game (Games [2016]). Overcooked is a multiplayer cooperative game where the goal is to work in a kitchen as a team with partner cooks and prepare various dishes together within a limited time. However, the game is highly dynamic. In many maps, the kitchen itself is not static and can change during a run. Random events such as pots catching fire add to the chaos. The challenge is to coordinate with the rest of the team and divide up subtasks efficiently.

The aforementioned game was simplified and reimplemented in a simpler environment (Carroll) to serve a purpose of scientific common ground for studying multi-agent cooperation in somehow complex settings. Many additional features of the original game were removed, leaving only essential coordination aspects. In its simplest form, the environment takes place in a small static kitchen layout, where the only available recipe is onion soup, which can be prepared by putting three onions into a pot and waiting for a given period of time. Somewhere in the kitchen there is an unlimited source of onions and a dish dispenser where the player can grab a dish to carry cooked onion soup to the counter. The team of cooks is rewarded as a team by an abstract reward of value 20 every time cooked soup is delivered to the counter. It may seem like a simple task. However, players face problems on several levels.

### 4.2 Basic layouts

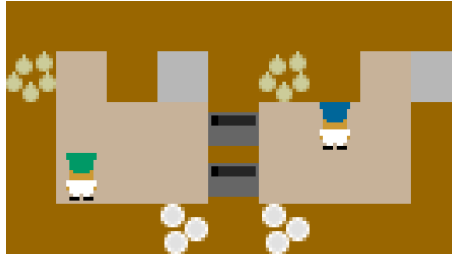
Although the Overcooked implementation has its own generator that can be used to generate new random kitchen layouts, most of the related scientific work to date has experimented with a fixed set of predefined layouts, each of which captures some important aspect of coordination.

#### Cramped Room



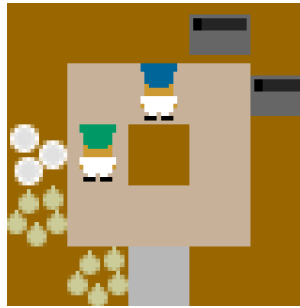
Cramped Room, as the name suggests, represents a cramped kitchen layout where all the important places are relatively easy to reach. The challenge lies in the low level of coordination of movement with the other partner, as there is no free space around.

## Assymetric advantages



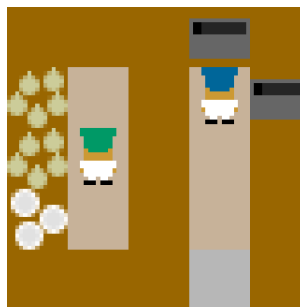
In Assymetric Advantage, both players are in separate regions where each region is completely self-sufficient. However, each region has better potential for a specific subtask. And only when both players make the best use of their own region's potential, the maximum joint efficiency is achieved.

## Coordination Ring



The Coordination Ring is another example of a layout that requires skillful coordination, as the only possible movement around the kitchen is along a narrow circular path that can be used in only one direction. For example, if one player chooses to move in a clockwise direction, the other player would automatically get stuck if they tried to move in a counterclockwise direction.

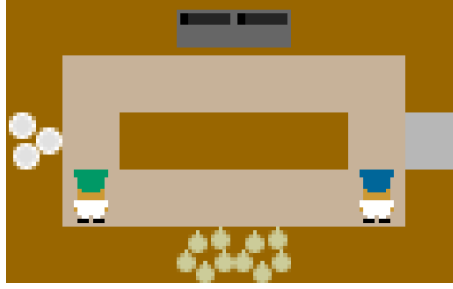
## Forced Coordination



The Forced Coordination kitchen layout is quite different from the others. In this layout, each player is located in an isolated region where no player has all the resources necessary to prepare a complete onion soup alone. Thus, players are forced to cooperate with each other using the resources they have at their disposal.



## Counter Circuit



In the last layout, the situation may look similar to the Coordination Ring. In this case, however, carrying onions around the entire kitchen is highly suboptimal no matter which direction the players choose. To deliver onions efficiently, players must pass them over the counter to shorten the distance. However, the cooks still need to decide who is responsible for bringing the plates.

## 4.3 Environment description

### 4.3.1 Action space

The action space is quite trivial, since it contains only the essential actions needed to operate in this environment.

- Go north
- Go south
- Go east
- Go west
- Stay
- Interact

### 4.3.2 State representation

There are two state representation functions predefined by the authors, namely `featurize_state` and `lossless_state_encoding`.

The first of the two extracts manually designed features into a one-dimensional vector of ones and zeros. These features can be interpreted as a partial observable representation of the environment state, as the majority of features are computed in relation to the nearest point of interest. For example, only the closest source locations of onions and dishes related to the player's current location are included, resulting in the loss of global information about the entire layout.

The second state representation, as the name suggests, provides lossless global information about the state of the environment. As in the previous case, environment variables are also represented by ones and zeros. However, unlike the previous function, the resulting representation is not a one-dimensional vector.

The result is formed by stacked masks, where each mask represents some feature of the environment in the form of a two-dimensional vector corresponding to the layout width and height. For instance, the mask representing the location of player  $i$  is a vector of shape  $(width, height)$  filled with zeros at all positions except for the value one at the coordinates where player  $i$  is located. Similarly, the mask of the same shape representing onion sources is filled with zeros, and at all coordinates where there is an onion source in the environment, there are ones instead. We can see that this representation goes beyond the nearest locations and provides global information.

### 4.3.3 Rewards

As we said in the introduction, the environment is purely cooperative in the sense that players share the common reward of value 20 each time a soup is delivered to the counter location. And this reward is player independent. The cycle of the environment is virtually infinite, since there is no a priori end state. Thus, it would be theoretically easy to get a cumulative sum reward of infinity. This is prevented by setting a finite time horizon, which by convention is set to the limit of 400 steps. With this constraint, it makes sense to compare the results of different runs.

In addition to the general reward for delivering soup, there may be player-dependent partial rewards that are not accounted for in the total sum of rewards. These rewards are mainly used for learning purposes of the agents, since especially on some particular maps it is highly unlikely that by following the initially random policy the whole process of soup making including delivery will be completed. Using predefined partial rewards

- Dish pickup reward - Dish is picked in useful situation, e.g. pot is ready or cooking.
- Soup pickup reward
- Placement in pot reward - Useful ingredient is added to the pot, in our setting only onion soup recipe is used, so this corresponds to the action where onion was added to the pot.

allows agents to learn subtasks first. It is important to eventually ignore these partial rewards during training, as agents could focus on these subtasks and ignore the main goal thus failing the main task altogether. This is usually implemented by linearly decreasing the weight of partial rewards over some finite time horizon.

### 4.3.4 Initial state

As there are two functions predefined by the authors to represent the state of the environment, there are also two ways to create the initial state of the environment.

The first of the two uses fixed layout initial player locations, which always creates an identical initial state including both players' initial locations. However, learning a single agent by always placing it in the same initial location will likely cause it to fail when starting in the opposite location. For this reason, *player\_index* is introduced to specify which agent is interpreted as player number

one and player number two. This index is randomized every time the environment is reset, ensuring that the agent encounters both starting locations during learning. Using this index introduces a minor chaos into the environment, as various parts of the environment control have to address the player index problem, and both actions and state observations have to be switched to the correct order, making it a bit opaque. We assume that this kind of initial state is suitable for some kind of benchmarking where the initial conditions are always the same.

In our experiments, however, we use the second approach to state initialization, where initial positions are always randomly sampled for both players. We claim that this is consistent with our intuition, since robustly cooperative agents should be able to cooperate well regardless of their initial position, which has been partially demonstrated (Knott et al. [2021]). However, this may come at the cost of reduced average performance. Besides random positions, the randomized initialization function also provides a probability threshold argument that can be used to randomly initialize some additional random objects in the environment.

## 5. Related work

A common approach to two-player games is to train an agent by confronting it with a set of other AI agents. This has been shown to perform impressively well against human experts in several complex games such as Dota (OpenAI [2019]) or Starcraft (Deepmind [2019]). The authors of the overcooked environment (Carroll et al. [2020]) believe that the distributional shift from AI training to human evaluation is successful due to the fact that the nature of this environment is strongly competitive. This is illustrated by the canonical case of a two-player zero-sum game. If humans take a branch in the search tree that is inadvertently suboptimal while in the role of minimizer, this only increases the outcome for the maximizer.

However, the situation is not so ideal in the case of common payoff games such as overcooked environment. If a self-play trained AI agent is paired with another sub-optimal partner, the result can be abysmal. In this case, both agents play a maximizing role in the search tree. If the self-play AI agent expects its sub-optimal partner to be the same agent, it can choose branches that lead to maximum joint payoff. However, since the sub-optimal partner may behave sub-optimally, it may happen that it mistakenly chooses some worse branches that are unforeseen by the self-play agent, which may lead to a significantly worse result. However, since the common payoff is shared between the two agents, this is no longer an advantage for the self-play maximizing agent. Rather, it leads to failure for both agents.

### 5.1 Human cooperation

Most previous work in this area has focused on one of two types of coordination. The first is coordination between a human and an AI partner. And second, focusing solely on the fully AI-driven pair.

While perfect AI-human coordination is generally a more desirable goal to achieve in all sorts of domains, it will not be our main focus. Several previous scientific papers have addressed this issue. A particularly noteworthy contribution is the paper On the Utility of Learning about Humans for Human-AI Coordination (Carroll et al. [2020]). They collected several human-human episodes and incorporated these experiences into training. This human data was used to create human-based models using behavior cloning and Generative Adversarial Imitation Learning (Ho and Ermon [2016]) and then incorporated into training. One of the important conclusions was that when self-play and population-based agents were paired with human models, the overall results were substantially worse than when paired with agents designed to play with human models.

### 5.2 AI-AI cooperation

We have already mentioned the problem when a self-play agent is paired with a suboptimal (e.g. human) partner. However, a similar problem can also occur when two different agents trained in self-play mode are paired together. A com-

mon approach used in competitive settings is to introduce diversity by exposing the trained agent to a diverse set of partners during the training phase. There are several popular methods for partner selection during training.

### 5.2.1 Self-play

The self-play method has already been mentioned, and it is arguably the simplest and most straightforward method for partner sampling. As the name indicates, the idea of this approach is that the agent currently being learned is paired with the copy of itself during training. This method doesn't introduce any kind of diversity into the learning process, as it only learns from its own current policy.

### 5.2.2 Partner Sampling

Partner sampling is an extension of the previous method. However, instead of playing only with the current policy, the partner is sampled from previously periodically saved policies of this given training period. The diversity in this method comes from the theoretical point of view that previous policies correspond to different behaviors.

### 5.2.3 Population-Based Training

The population-based training method is based on an evolutionary algorithm that focuses on hyperparameter training and model selection. The population consists of agents parameterized by a neural network and trained by a DRL algorithm. During each iteration, agents are drawn from the population and trained using the collected trajectories. Pair-wise results are recorded, and at the end of the iteration, the worst agents are replaced by a copy of the best agents with their hyperparameters mutated.

### 5.2.4 Pre-trained Partners

Finally, a often employed method is to employ pre-trained partners. Diversity is expected here due to the fact that different runs of reinforcement learning algorithms often yield different agent behaviors.

### 5.2.5 Result

The problem with fully cooperative games is that when agents are trained together, they tend to exploit the common knowledge they have learned during training, which often makes them unable to cooperate with unseen agents. It has been shown (Charakorn et al. [2020]) that this is true for the first three methods mentioned, where, when performing a cross-play evaluation of the set of different agents obtained by the same method, only those pairs of agents that were explicitly trained together performed well, while the remaining pairs failed.

Nemam si tady od nich z toho clanku pujcit primo ty obrazky ukazujici ten SP off-diagonal failure, abych se na to pak mohl vizualne odkazat ve svých experimentech? Nebo krast obrazky se nedela? The authors conclude with results

showing that incorporating different pre-trained agents for training robust agents was significantly more successful than other methods.

## 5.3 Problem of robustness

Another important issue related to common payoff is evaluating the performance of trained agents. As discussed earlier, we cannot rely on training performance alone, no matter what metric is chosen, because training sets of agents exploit shared knowledge and thus make training performance appear excellent. As noble as it sounds, it is quite complicated to define which agent behavior is actually robust. Obviously, it's a desirable goal to create a trained agent that can cooperate well with all possible behaviors of its partner. However, since it is difficult to come up with a diverse set of agents for partner sampling during training, the same is true for evaluation agents.

### 5.3.1 Average performance

The first obvious choice of evaluation metric could be to look at the average evaluation performance over the evaluation set of agents. It is questionable, though, whether this answers the question of robustness. In one scenario, we could theoretically achieve a result where our agent fully fails to cooperate with half of the evaluation agents, while performing excellently with the other half. In another case, we could achieve some kind of intermediate performance with the whole set of evaluators, where cooperation does not fail completely, but performance is also nowhere near the optimal values. In both cases, the average evaluation performance reaches similar values.

### 5.3.2 Threshold performance

Another possible approach for the given layout could be to define our own threshold, which could be practically equivalent to a certain number of delivered soups. The evaluation metric could then simply be expressed in terms of the number of evaluation agents that managed to reach such a threshold when paired with our trained agent.

### 5.3.3 Edge case testing

The problem with previous metrics is that they are reasonable only in a situation where the agent's evaluation set is complex enough in terms of behavioral diversity, which also implies a condition on the sufficient size of such a set. However, creating an evaluation set that satisfies this condition is practically impossible.

Instead of trying to construct such an evaluation set and trying to evaluate whole episodes in an effort to test cooperation on the run, we could break the interesting cooperation-challenging situation into a number of separate tests (Knott et al. [2021]). Suddenly, our thinking about robustness shifts from looking at cooperation from the point of view of the outcome of the entire episode to just separate small modular situations where the desired correct behavior of both agents can be described. This approach is inspired by unit testing in software

development, where even if the program behaves correctly in the vast majority of cases, it can still produce undesirable responses in edge cases.

Similar to unit testing, we can define a wide set of situations (e.g., instances of overcooked environment states) and pair them with a set of acceptable consequent behaviors of agents that can be labeled as robust response. Evaluation using such a metric can then be expressed as the number of tests that have passed.

However, there are several challenges with this metric as these edge cases require manual design, in many scenarios it is also difficult to know with certainty what the correct expected behavior is, and finally, similar to unit testing in software development, there are never enough defined edge cases to cover all possible situations.

## 6. Our work - Preparation

Before we get into our own experiments with training a robust agent, which will be the main focus of the next chapter, we want to prepare our solution first and reimplement some of the previously mentioned results regarding self-play agents.

### 6.1 Framework

Writing all of our code on top of the overcooked environment from scratch is not necessary, as there are already several frameworks that implement deep reinforcement learning algorithms. Although in general there are many more DRL frameworks that could be suitable for our purposes, most related projects using the overcooked environment have so far used either the RLib library or Stable Baselines.

As the authors themselves say, RLib provides support for production-level and highly distributed RL workloads. This framework is also suitable for multi-agent learning. The downside here is that it is less suitable for smaller projects and development on a local machine, as there is a fair amount of overhead due to its parallelization capabilities towards cluster computing. Personally, I found the framework a bit intimidating and the documentation a bit confusing.

Another option was the Stable Baselines framework, which in my opinion offers documentation that is clearer and easier to understand. Also, in my humble opinion, the code base structure is more transparent and basic API is quite light. To demonstrate its simplicity, once you have an environment that implements the standard RL OpenAI environment interface, you can perform all training with default hyperparameters and policy represented by a multi-layer perceptron network using the PPO algorithm as simple as follows:

```
env = make_vec_env("CartPole-v1", n_envs=4)
model = PPO("MlpPolicy", env, verbose=1)
model.learn(total_timesteps=25000)
```

The downside here is that the framework does not have native support for simultaneous multi-agent learning.

#### 6.1.1 Modifications

After considering both frameworks and the scope of our project, we decided to build our project within the Stable Baselines framework. We will not cover all the details of our diverse population approach, as there is plenty of room for this in the next chapter. Here, we just want to mention some of the major framework modifications that had to be carried out in order to bend the framework to suit our purposes.

##### Embedding partner

First, as mentioned above, the Stable Baselines framework does not support multi-agent learning. In our approach, we will make a slight simplification by theoretically transforming the environment from a multi-agent to a single-agent



setting. Of course, we will not reduce the number of agents in the environment, as it will still be multi-agent. Rather, we will look at the situation from the point of view of the single agent that is learning. In this way, we can look at the environment as if the partner cook was only embedded as part of the system.

And exactly the same change is necessary in the framework. During a training run, the partner is embedded as part of the environment, and at each step of the loop, the actions of the partner are sampled. The actions of the trained agent are sampled in the usual way, as expected, while the actions of the partner cook are sampled according to the policy of the embedded partner. The resulting two actions are concatenated and passed to the environment as expected. This mechanism allows us to use any partner sampling method (5.2), including self-play, where the same instance of the parameterized policy is used as the one being learned.

### **Convolution policy**

By default, Stable Baselines provides several types of parameterization policy representations, including fully connected dense multi-layer perceptron (MLP) and convolutional neural network (CNN). However, regarding the CNN representation, when using this wrapper, the framework expects the inputs to be strictly in a standard image format (RGB, RGBD, GrayScale). Unfortunately, this is accompanied by the use of inadequate assumptions throughout the code base, heuristically expecting the inputs to be in a certain format. For example, the function to detect if the image space is in the first format only checks if the first dimension is the smallest. This erroneously returns true for some overcooked layouts using lossless state representation (4.3.2) where width is less than both height and definitely the number of stacked masks. In lossless encoding, masks are equivalent to channels and are represented by the last dimension. Authors might argue that with such a state representation, the inputs are technically not images, and therefore such sanity checks expecting images should not be used in the first place. However, this is the problematic part as I was unable to find out whether the framework allows some options to bypass these constraints. Eventually, after disabling these assertions in certain critical places, the code managed to work flawlessly with our convolutional representation without any further modifications.

### **Loss and rewards augmentation**

Finally, two important aspects of our experiments will revolve around augmenting rewards and extending the objective loss function. These modifications are not truly related to Stable Baselines specifically, as these would have to be made regardless of the framework chosen, we just want to mention them here to make the list of changes complete. Fortunately, these modifications were also quite straightforward, as both of these elements are nicely separated in transparent, easily extensible locations.

## 6.2 Self-play

Before we dive into our experiments, we want to make sure that our solution is ready to be used with the environment. As a reasonable setup, we'll start by reimplementing basic agent training using self-play (5.2.1) and try to demonstrate its allegedly poor cooperation abilities.

### Default setting

We started with the simplest default settings. For initial experiments we chose to start with the first layout Cramped Room (4.2) using the more lightweight of the two, partially observable state representation (4.3.2). As we recall, this representation uses a one-dimensional vector of features that is out-of-the-box compatible with prefabricated Stable Baselines MLPPolicy wrapper, which contains two shared hidden layers before separate layers for actor and critic are applied. Despite several attempts to tweak the hyperparameters, we were not able to come close to the reported results of cooperation failure (5.2.5). In all of our experiments, different self-playing agents managed to cooperate reasonably well. While there were some pairs of agents that did not cooperate at all, this was rather rare compared to the overall cross-play evaluation (Figure 6.1). However, we believe that the problem of uncooperative behavior may be exacerbated by this evaluation scheme in some layouts more than others, since different layouts present different challenges to cooperative behavior. To make matters worse, while we were able to obtain similar cooperative behavior in the Asymmetric Advantage and Coordination Ring layouts, the same training procedure failed completely in the remaining Forced Coordination and Counter Circuit layouts. Here, the self-play training procedure was unable to learn how to perform the main task at all.

### Lossless state representation

To both overcome this problem and make our experiments as relatable as can be to previous related work, we tried to adapt our solution as closely as possible to the settings used in previous work. We modified our solution based on the original overcooked environment project (Carroll et al. [2020]), whose specific settings are described in detail in Appendix A. Instead of the partially observable state representation, the global lossless representation was adopted. This implied the need to change the policy parameterization as well. The Stable Baselines framework provides a prefabricated convolutional CNN policy wrapper. However, the convolutional part of the network had to be manually modified to match the described structure. [mam se poustet do detailniho popisu zmen, nebo staci ze jsem zminil appendix A, kde to lze dohledat?](#) In addition, several PPO hyperparameters have been changed compared to the default stable baseline PPO setting.

### Hyperparameters search

With these settings adjusted, we attempted to run the self-play agent training on all layouts as before. On the basic Cramped room layout, the training seemed

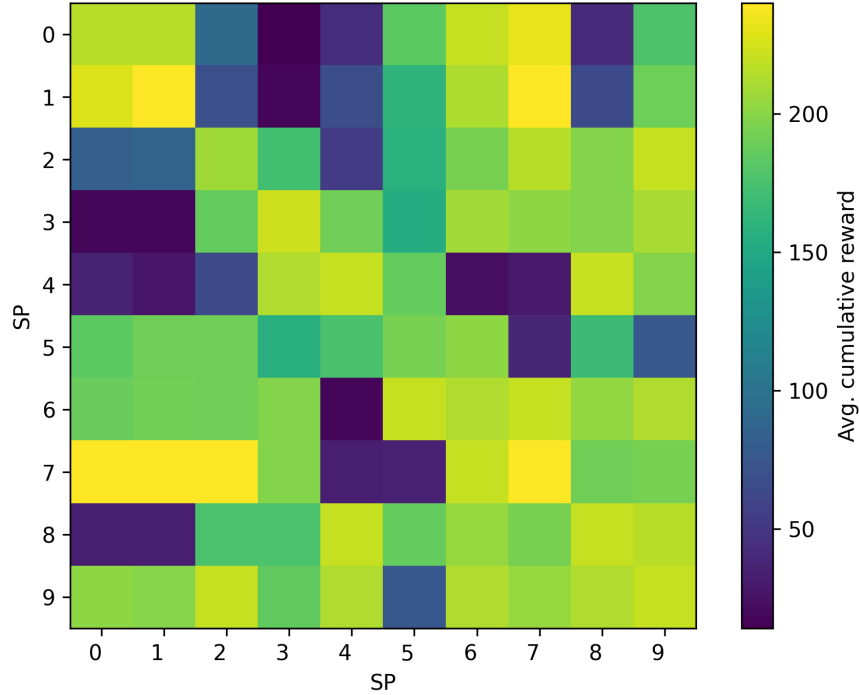


Figure 6.1: Self-play MLP cross play evaluation

Cross-play evaluation of 10 agents trained via self-play method on cramped\_room layout, where policy is parametrized by MLP

to progress successfully. However, when cross-play was evaluated, similar pairwise cooperative behavior was achieved as in the previous state representation. However, training failed on all other layouts.

We performed hyperparameter search by running the training process with several hundred randomly different hyperparameters to find a suitable stable configuration. With the resulting configuration (table 6.1) we were able to train a self-play agent on all layouts in most of the runs.

## Preemptive divergent check

Unfortunately, even this configuration is not completely stable. In some runs, training on more complex layouts such as Forced Coordination and Counter Circuit fails to achieve the main goal. In these cases, the agents always learn how to successfully acquire partial rewards by performing the subtasks, but fail to learn the final soup delivery. We suspect that the combination of layout complexity and likely policy exploitation prevents the agent from experiencing the last part of the task as often as necessary to adapt such behavior. If this is the case, it is highly unlikely that the agent will discover the soup delivery behavior once the time steps reach the partial reward shaping horizon. To account for this fact, we perform a preemptive divergence check at the  $3 \cdot 10^6$  time step, where we measure the average number of soups delivered during training. We have heuristically found that the average soup reward threshold of 3 over the collected episodes is sufficient to ensure that the agent will learn the main goal in the remaining training process.

<b>Name</b>	<b>Originally proposed</b>	<b>Applied</b>
<i>PPO Hyperparameters</i>		
Discount factor $\gamma$	0.99	0.98
GAE factor $\lambda$	0.98	0.95
Learning rate	$10^{-3}$	$4 \cdot 10^{-4}$
VF coefficient	0.5	0.1
PPO clipping	0.05	0.1
Maximum gradient norm	0.1	0.3
Gradient steps per minibatch	8	8
Minibatch size	2000	2000
Number of parallel environments	30	30
Total timesteps	$6 \cdot 10^6$	$5.5 \cdot 10^6$
Entropy bonus start coefficient	0.1	0.1
Entropy bonus end coefficient	0.1	0.03
Entropy bonus horizon	none	$1.5 \cdot 10^6$
Preemptive divergent time step	$3 \cdot 10^6$	$3 \cdot 10^6$
<i>Environment settings</i>		
Episode length	400	400
Soup delivery shared reward	20	20
Dish pickup partial reward	3	3
Onion pot placement partial reward	3	3
Soup pickup partial reward	5	5
Partial reward shaping horizon	$2.5 \cdot 10^6$	$2.5 \cdot 10^6$
Initial state player locations	static	random

Table 6.1: Hyperparameters for self-play agent training on Cramped room layout

### 6.2.1 Random state initialization

Despite these adjustments, there remained problems on two particular layouts. On the Asymmetric Advantages layout, the resulting performance of the trained agent was significantly worse than in the presented results (Carroll et al. [2020]). Whereas on the Forced Coordination layout, training converged to zero performance suspiciously often. We found that there was a missing constraint on the sampled initial positions of the agents, which allowed the initial positions to be in the same isolated layout region. While this is not a problem for the remaining layouts, it is an important obstacle for the aforementioned maps.

In both Asymmetric Advantages and Forced Coordination layouts, not all locations are globally reachable. As a result, in Assymmetric Advantages, the agent was forced to learn how to operate within its own region, rather than relying on beneficial cooperation from its partner using its region advantage. And this led to a significantly lower performance value.

In the Forced Coordination layout, the situation was even worse, because if both players started in the same isolated region, they would not be able to finish the soup at all.

Lastly, the authors (Knott et al. [2021]) believe that introducing random initial states can lead to better robustness, which could be an explanation for our better measured cross-play cooperation. Nevertheless, we ran the same training with fixed initial positions and the results were moreless the same.

### 6.2.2 Focus on Forced coordination

With all the observations we have made thus far, we can state that there are vast differences between the different layouts. Therefore, for our subsequent experiments, we are primarily limiting our attention to one particular layout. We believe that the Forced Coordination layout provides a good balance between complexity and also emphasizes the importance of cooperation, which is the main point of our interest. Furthermore, this layout shows signs of a cooperation problem that resembles the reported result (5.2.5) when different self-play agents are evaluated using cross-play (6.2).

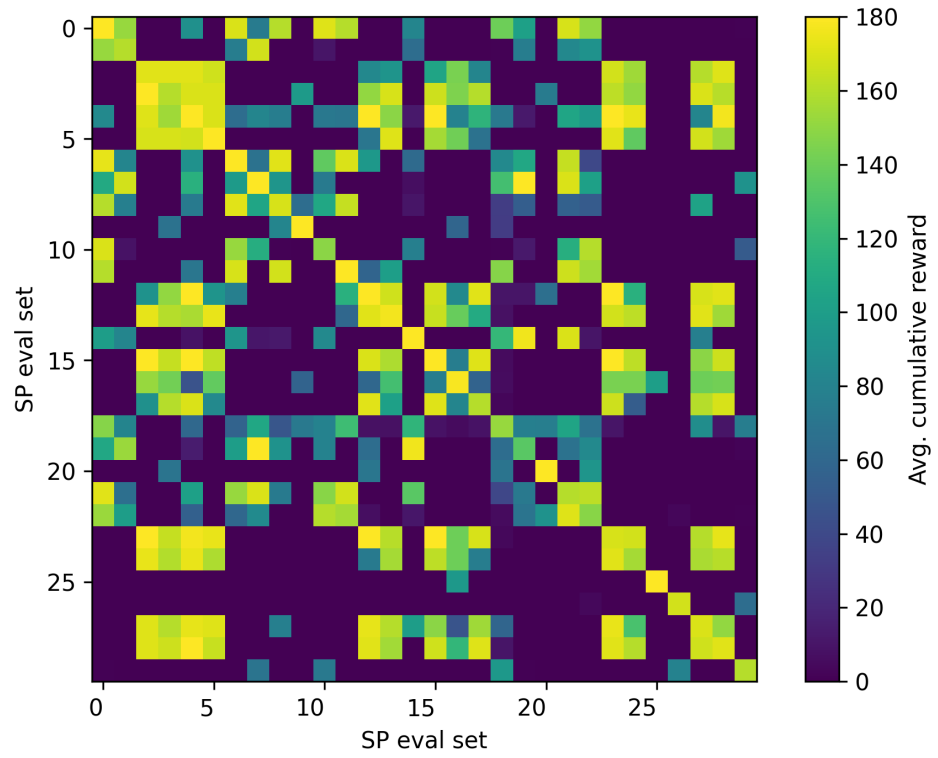


Figure 6.2: Self-play CNN cross play evaluation on Forced coordination

Cross-play evaluation of 30 agents trained via self-play method on Forced coordination layout, where policy is parametrized by convolutional network

## 7. Our work - Contribution

In this chapter, we propose a novel approach to injecting explicit diversification into agent behavior through diverse population training. Previously mentioned traditional population-based methods were not effective enough to add the necessary diversification to the population. We try to use population training in slight different manner with our primary focus set on two important aspects. Firstly we try to face well known problems connected with multi agent system by changing our perspective towards single-agent settings and transforming to iterative process. And secondly we try to address the lack of explicit diversification in standard population-based learning methods by pursuing population policy diversification.

In following sections we attempt to walk through several different aspects of our population building approach. These parts can be seen as a soft guideline towards diverse population training algorithm which will be presented at the end of this chapter.

### 7.1 Multi agent setting simplification

#### 7.1.1 Single-agent perspective

As we could have already seen (3) multi agent systems can be much more intimidating to tackle down. One of the proposed suggestion was to try to simplify the sitatio by looking at the multi agent environment from the single-agent point of view where the goal is to learn joint policy. As discussed before, this is often impossible and mostly inpractical.

Nevertheless we try to, in a sense, attempt to utilize similarly motivated approach. We propose single-agent learning approach where at any given time only one agent is being learned. Naturally we cannot simply ommitt the second player from the environment as we cannot deny the multi agent essence from the environment. We rather look at the environment as being instance of single-agent environment where the parnter is embedded and in a sense forms an integral part of the system.

#### 7.1.2 Non-stationarity

We propose population training to be an iterative process where only one agent is being learned and only after its's training is finished he is included into the population. Once an agent is added into the population it's policy is fixed and no further updates are applied for the rest of population training process. Only agents used in environment embedding are the ones coming from the already trained population of agents. This claim this effectively solves the non-stationarity problem we encountered in multi agent settings as from the point of learning agent the behavior of its partners never changes.

There is unfortunately downside to this approach. While previously the environment had its state transition function purely deterministic, now with the partner embedded into the system it becomes stochastic. Nonetheless, we hypoth-

esize that this may also be in favor of the final robustness, since by experiencing stochastic state transitions, the agent could theoretically be more robust as a result, since it will have to learn how to behave in a more stochastic exploratory environment.

### **7.1.3 Population exploitance**

We hypothesize that by utilizing iterative process with past agents being fixed it could also help with population exploitance. When set of agents are trained concurrently together they may all exploit their behavior towards shared experience. It could probably happen that the new agents might still try to overfit their behavior towards the population agents. However, by fixing previous past agents we try to break this relation at least from one side of the couple. We are optimistic that this, in conjunction with further diversification factors, will greatly reduce exploitation of the population.

## **7.2 Population building**

A mere change in our perception towards single-agent settings would probably not be enough of an improvement towards agent robustness. Therefore we propose several ideas to promote population diversification.

### **7.2.1 Inicialization**

To start with, inspired by original motivation behind all population-based methods, we follow the idea that populations should be in general always more diverse than single agent. Hence when utilized in training this diverse robustness should be passed onto the next learned agent. In order for this to be true the previous condition on non existence of population exploitance must hold. Additionally we already proposed iterative training process where we suppose preexistence of some set of agents utilized for partner embeddedding.

We propose that in order for other diversification factors to work we have to start with some initial population for partner embedding sampling that is already somewhat diverse. One could argue that if we had a method for creating such initial diverse population we could just stop there and extend this population using such method. However, we do not see the perfect diversity and robustness of initial population to be crucial. We rather expect the inital population to introduce somewhat mildly different set of behaviors just to give a resonable starting point for the rest of population training.

### **7.2.2 KL divergence**

We hypothesize that preexistence of diverse population is just a soft precondition. We propose that the main ingredient behind our diversity building process lies in attempting to shift the behavior of each consequentially trained agent to be as different as possible from the behaviors in the existing population. Since sampled population partners are embedded into the environment we can look at this approach as domain randomization technique.



To emphasize the difference between two behaviors, we must first be able to measure it. There are arguably various possible ways how to express difference between two behaviors including methods discussed in evaluation approaches (5.3). Nevertheless since the agent behavior is represented by parametrized policy we suppose that one of the more straightforward methods for low-level difference could be to look at the difference between the two policy distributions.

When measuring the distance of the probability distribution the Kullback-Leibler divergence (KLD) often comes in mind.

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

Although it is not a metric since it is not symmetric in the two distributions and also it does not satisfy the triangle inequality, we can still utilize it as we are mostly interested in relative distribution differences rather than absolute values. We propose idea of incorporating maximization of KL divergence as a tool to differentiate new agents from the ones that are already contained in the population by taking the average of KL divergence of trained policy and policies from population.

We suggest that there are two different suitable ways how to incorporate KL divergence maximization into the process. Firstly, KL divergence between policy that is currently being trained and the policies in the population can be used during episode sampling by interpreting it as another additional partial reward based on a current state of the environment. And secondly, similar approach can be applied on the level of PPO objective function where another additional objective term based on maximization of KL divergence can be added.

It may seem that these two approaches might have the same effect. However, we believe that they work in a slightly different manner. When applied in the PPO objective the term is trying to differentiate current policy no matter the current state of the environment. While applying the KL divergence bonus as partial reward based on the state might work as a exploration enhancing technique since it should push the agent to visit more the states where it's current policy differs the most which could potentially lead to state trajectories that have not been seen by the agents in the population.

## KL divergence coefficients

Applying the absolute value of KL divergence the way it is calculated using the formula would not work in our favor as we have to consider the value ranges with regard to the environment properties. KL divergence is non negative but it can easily reach values of infinity when not bounded. Additionally we need to consider suitable multiplicative coefficient to scale the values to reasonable range.

When new parametrized policy  $p$  is initialized it produces the probability distribution that is close to uniform. Having sampled several already trained self-play agents and evaluated in same initial states we obtained specialized policies yielding the probability distributions  $q_1, q_2, q_3$  (see table 7.1).

These values of KL divergence can give us an idea about the absolute values when new policies are compared to already trained agents. We believe that probability distributions  $p$  a  $q_1$  are different enough in terms of what we are

	Probability distribution	$D_{KL}(p  q)$
$p$	[0.1679, 0.1647, 0.1644, 0.1663, 0.1713, 0.1654]	
$q_1$	[0.0336, 0.0729, 0.0467, 0.5419, 0.1431, 0.1617]	0.449
$q_2$	[0.0453, 0.1906, 0.1270, 0.2221, 0.0942, 0.3207]	0.183
$q_3$	[0.1875, 0.1851, 0.1663, 0.1369, 0.1414, 0.1828]	0.009

Table 7.1: KL divergence values illustrated on trained policies on same random initial state

trying to achieve when training population of diverse agents. By considering the absolute value of KL divergence of these two distributions we try to derive suitable KL divergence coefficient for loss and reward augmentation.

We hypothesize that KL divergence term which we add to PPO objective should have comparable weight as the entropy bonus (recall 2.5) as this bonus is also designed only to be mild exploratory pressure rather than main objective. Entropy in early stage of training corresponds to uniform probability distribution which yields absolute values of 1.8. This combined with entropy bonus start coefficient of 0.1 (hypereparameters table 6.1) produces absolute values of cca 0.18. At the final stage when policy is already trained the entropy usually decreases to absolute values of around 1.2 which when combined with entropy bonus end coefficient of 0.03 produces final value of 0.036.

As far as reward augmentation is considered we believe that when all KL divergence bonus are summed throughout the 400 steps of episode it should correspond only to fraction of what can be obtained by following main objective of soup deliverence. We propose to design such reward clipping in equivalence to the number of soups delivered (0.5, 1, 1.5 (R0, R1, R2, respectively)).

After performing several trivial experiments we concluded with following experiment settings (table 7.2) that we want to evaluate during our following experiments.

Experiment Name	Reward coefficient	Reward clip	Loss coefficient	Loss clip
No divers.	0.0	0.0	0.0	0.0
R0	0.08	0.025	0	0
R1	0.15	0.05	0	0
R2	0.1	0.075	0	0
L0	0	0	0.08	0.03
L1	0	0	0.12	0.07
L2	0	0	0.1	0.15
R0L0	0.08	0.02	0.08	0.02
R1L1	0.1	0.04	0.1	0.03

Table 7.2: Summary of experiments utilizing KL divergence

### 7.2.3 Population structure

Finally, we need to provide description of entire proposed population building process. We have already dealt with population initialization where we proposed how the initial agents could be constructed and we also designed diversification tools to enhance the behavioral diversity of the population.

However, our ultimate goal is to provide a procedure that produces final agent or set of final agents that should be robust and highly cooperative. By applying diversifying techniques we effectively train agents that are not fully optimizing towards main objective but their performance can be degraded due to diversification.

After having trained several diversified population agents we propose to train final agents by exposing them to already built population but omitting the diversification pressures. We propose to train first final agent by exposing him with population containing also the agents that were part of initialization and then train second agent using same population only with initial agent excluded.

To make sure the final agents have enough space to learn how to cooperate with entire population we extend the number of training timesteps to the factor of 1.5.

## 7.3 Simple convolution experiments

### Initial experiments visualised

After initial population training experiment (figure 7.1) we can see that first three self-play agents used for initialization of population play important role. In three seeds (10,11,13) initial population apparently strongly represents same one set of behaviors and lacks diversity towards other behaviors. As a result as agent training progresses during population building, agents learn to fixate solely on this kind of behavior as it promises maximal available outcome and ignore other types of behaviors.

To confront this issue we propose two suggestions. Firstly, in these initial experiments first three trained self-play agents were selected for population initialization automatically. We instead propose performing initial population agent selection manually by first training five self-play agents, evaluating them using any kind of evaluation method (in our case evaluating them with other set of self-play agents) and finally out of these five selecting such three agents that are both successful and also pair-wise different as much as possible.

And secondly, during training all available agents from current population were present at each training epoch as these were equally distributed among 30 parallel environments. This guarantees that the dominant type of behavior from population initialization was always present during training, which enabled newly trained agent to move towards this kind of behavior. Instead we propose not to sample all agents from current population, but rather sample only few of them to give chance to omit the dominant ones and allow trained agent to focus also on the remaining types. After performing simplified experiment we concluded with suitable value three for the number of sampled partners.

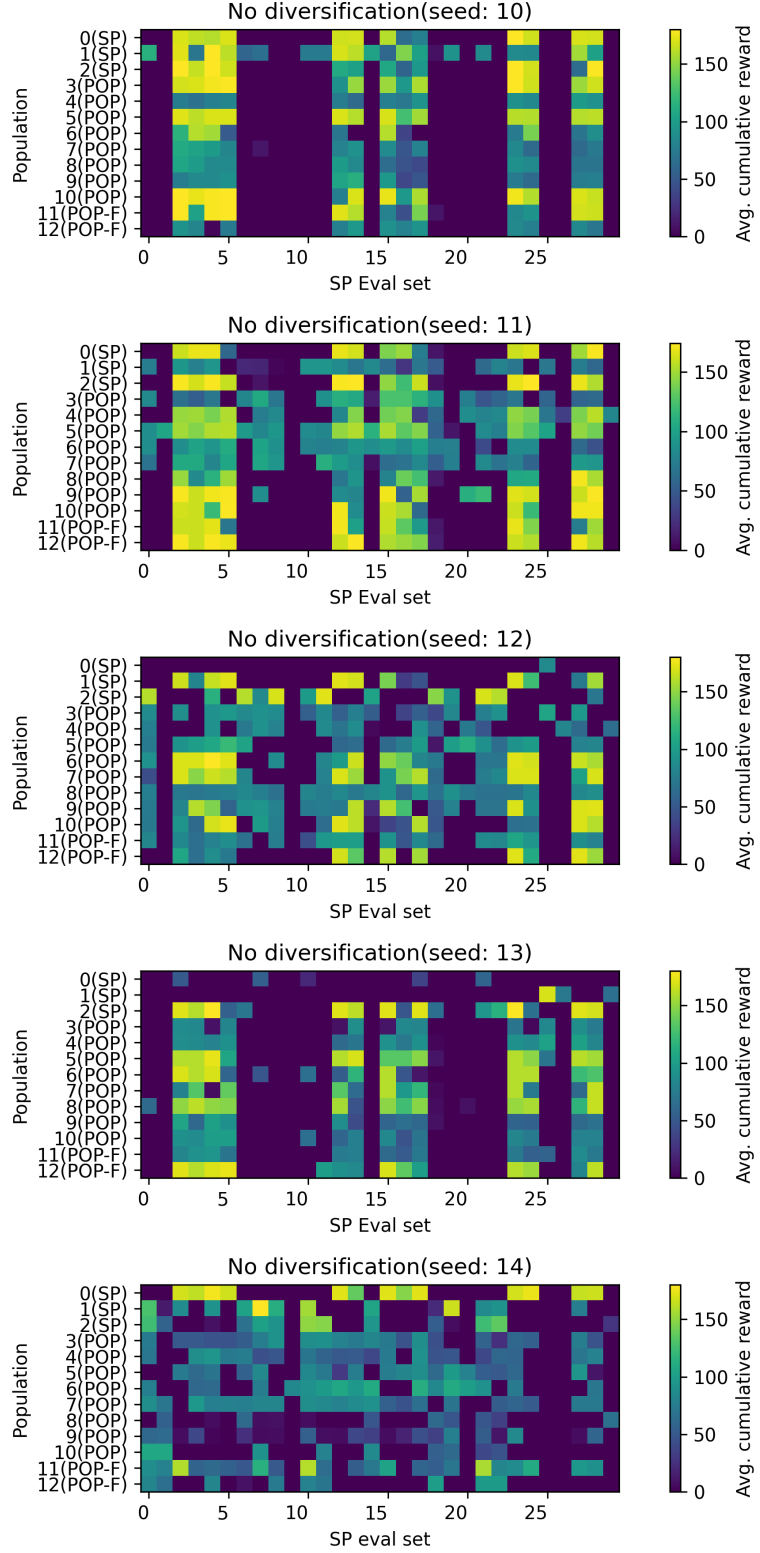


Figure 7.1: CNN Population with no diversification cross play evaluation  
 5 populations trained using different seeds with no additional diversification  
 cross-play evaluated against Self-play evaluation set

### Average cummulative reward evaluation

As discussed in previous chapters (5.3) performing evaluation of acquired population using usual methods is not very intuitive. When standard average cum-

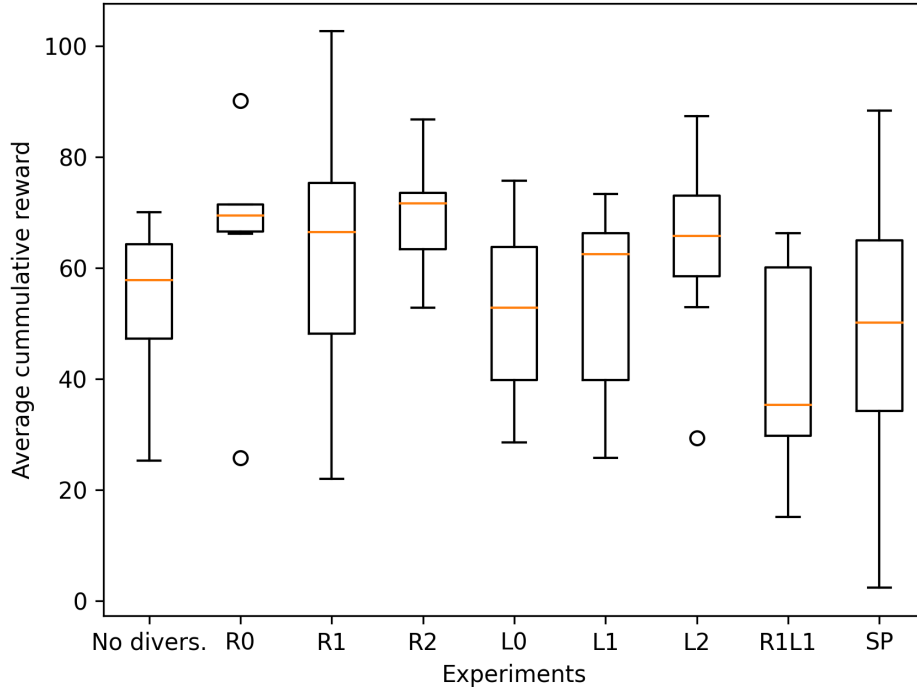


Figure 7.2: Average cummulative reward of final agents

Average cummulative reward of the two final agents trained in last phase of the population building. Includes values also for cross-play evaluation of self-play evaluation set for reference.

mulative reward metric was used (figure 7.2) there were just mild improvements achieved compared to self-play agents however results varied a lot which we believe could be attributed to two factors. Firstly, as we just discussed, the average cummulative rewards of initial populations also vary a lot when evaluated against self-play evaluation set. And secondly, even when final population agents are trained there is no guarante, that they will not end up in some suboptimal local optima with respect to current population.

### Ordered average cummulative reward evaluation

When agents are ordered by the evaluation with other agents it shows significant improvement in terms of early mild cooperation (7.3). Results are even more amplified when not looking at average case, but rather at 0.15 quantile (figure 7.4). *This might be interpreted as looking at worse case cooperation? In a sense that our final agent is able to cooperate to extent of cca 80 shared rewards with way wider set of evaluation agents?*

Sometimes even in the middle of population there can be already agents that behave nicely (7.6)

### Frame stacking experiments

TODO: two variants

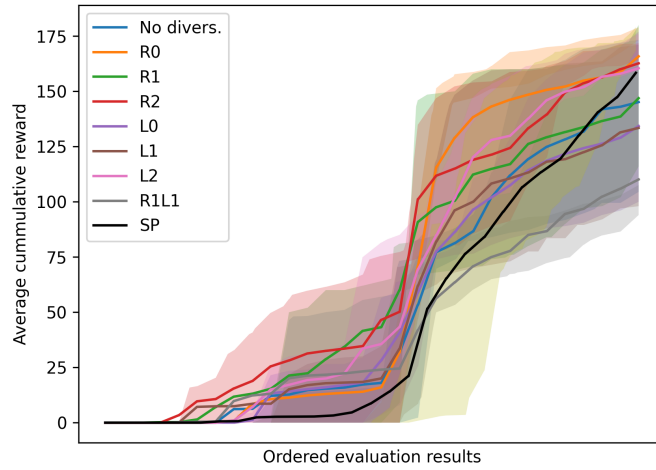


Figure 7.3: Average cummulative reward of final agents ordered by evaluation results

Average cummulative reward of the two final agents trained in last phase of the population building. Includes values also for cross-play evaluation of self-play evaluation set for reference. Average, 0.25Q, 0.75Q

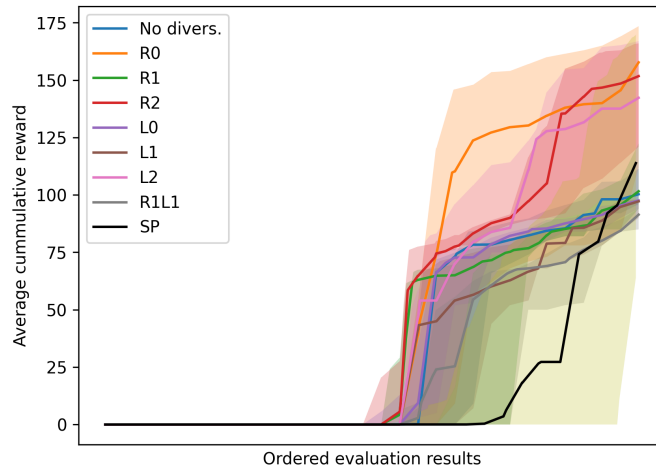


Figure 7.4: Average cumulative reward of final agents, 0.15quantile

Average cummulative reward of the two final agents trained in last phase of the population building. Includes values also for cross-play evaluation of self-play evaluation set for reference. 0.15Q, Min, 0.3Q

Looks like temporal feature made even basic self-play agent a bit more cooperative (figure 7.5).

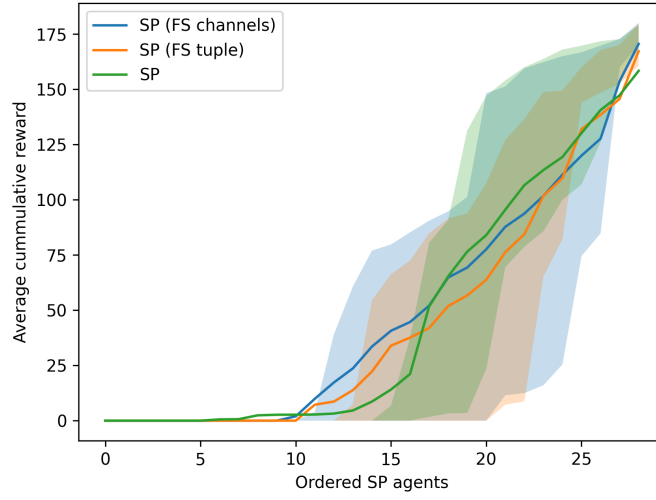


Figure 7.5: Average cummulative reward of final agents

Average cummulative reward of all agents from evaluation set. All three types of self-play agent designs included.  
Average, 0.25Q, 0.75Q

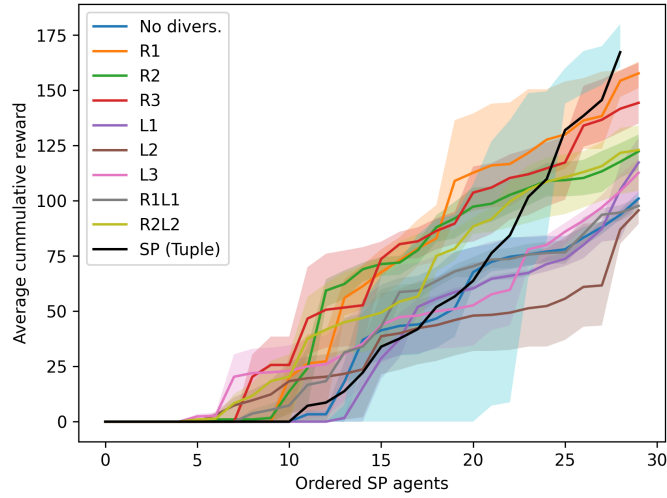


Figure 7.6: Average cummulative reward of 4th agent within the population part

Average cumulative reward of 4th agent from population. Using tuple frame stacking mechanism. Sometime the final pop agent is not the best one.  
Average, 0.25Q, 0.75Q

### 7.3.1 Populations evaluation

In here we will not evaluate agents against self-play agents, but evaluate final population agents against each other among different experiments settings. (Simple CNN 7.7). (FS Channels 7.8). (FS Tuple 7.9)

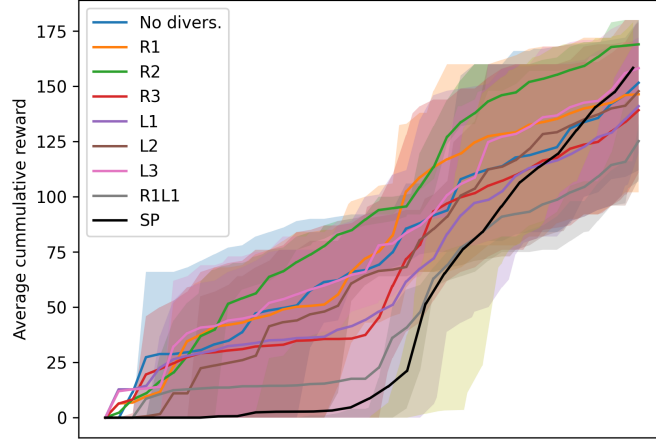


Figure 7.7: Evaluation of final population agent using simple CNN  
Average, 0.25Q, 0.75Q

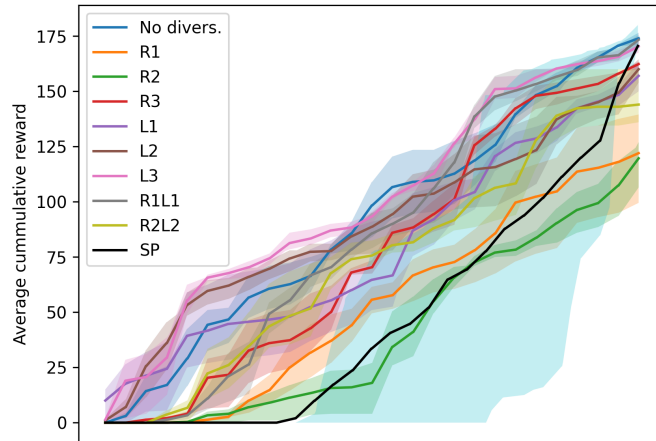


Figure 7.8: Evaluation of final population agent using channels frame stacking  
Average, 0.25Q, 0.75Q



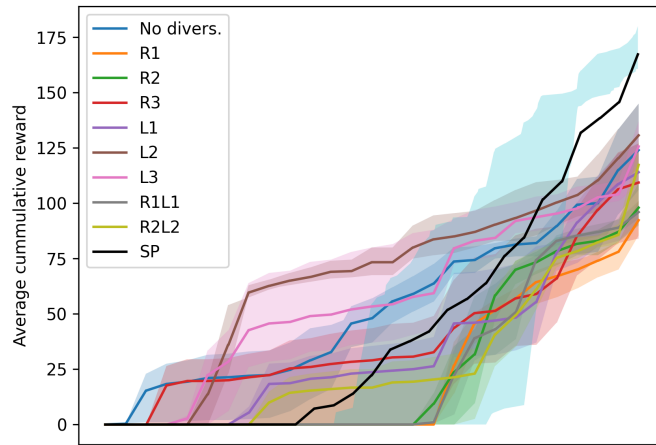


Figure 7.9: Evaluation of final population agent using tuple frame stacking  
Average, 0.25Q, 0.75Q

# Conclusion

So far it looks that our method works if we set threshold for cooperation to smaller value (cca 70 shared reward) and look for 0.30 quantil.

KL divergence as additional partial reward seems to be reasonable exploration enhancing technique that did not have that significant negative effect on trained policy. But when applied in objective function it seems to force the learned policy the way that is worsening the main objective.

# Bibliography

- Kullback leibler divergence. URL [https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler\\_divergence](https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler_divergence).
- Josh Achiam, 2018. URL <https://spinningup.openai.com/en/latest/>.
- Larry Armijo. Minimization of functions having Lipschitz continuous first partial derivatives. *Pacific Journal of Mathematics*, 16(1):1 – 3, 1966. doi: [pjm/1102995080](https://doi.org/10.1112/pjm/1102995080). URL <https://doi.org/>.
- Craig Boutilier. Planning, learning and coordination in multiagent decision processes. In *Proceedings of the 6th Conference on Theoretical Aspects of Rationality and Knowledge, TARK '96*, page 195–210, San Francisco, CA, USA, 1996. Morgan Kaufmann Publishers Inc. ISBN 1558604179.
- Micah Carroll. URL [https://github.com/HumanCompatibleAI/overcooked\\_ai](https://github.com/HumanCompatibleAI/overcooked_ai).
- Micah Carroll, Rohin Shah, Mark K. Ho, Thomas L. Griffiths, Sanjit A. Seshia, Pieter Abbeel, and Anca Dragan. On the utility of learning about humans for human-ai coordination, 2020.
- Rujikorn Charakorn, Poramate Manoonpong, and Nat Dilokthanakul. Investigating partner diversification methods in cooperative multi-agent deep reinforcement learning. In Haiqin Yang, Kitsuchart Pasupa, Andrew Chi-Sing Leung, James T. Kwok, Jonathan H. Chan, and Irwin King, editors, *Neural Information Processing*, pages 395–402, Cham, 2020. Springer International Publishing. ISBN 978-3-030-63823-8.
- Deepmind, 2019. URL <https://www.deepmind.com/blog/alphastar-mastering-the-real-time-strategy-game-starcraft-ii>.
- Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods, 2018. URL <https://arxiv.org/abs/1802.09477>.
- Ghost Town Games, 2016. URL <https://ghosttowngames.com/overcooked/>.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Soft actor-critic algorithms and applications, 2018. URL <https://arxiv.org/abs/1812.05905>.
- Matteo Hessel, Joseph Modayil, Hado van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning, 2017. URL <https://arxiv.org/abs/1710.02298>.
- Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *NIPS*, 2016.

- Sham M. Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, 2002.
- Paul Knott, Micah Carroll, Sam Devlin, Kamil Ciosek, Katja Hofmann, A. D. Dragan, and Rohin Shah. Evaluating the robustness of collaborative agents, 2021.
- Solomon Kullback. *Information Theory and Statistics*. Wiley, New York, 1959.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning, 2015. URL <https://arxiv.org/abs/1509.02971>.
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments, 2017. URL <https://arxiv.org/abs/1706.02275>.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei Rusu, Joel Veness, Marc Bellemare, Alex Graves, Martin Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dhharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518: 529–33, 02 2015. doi: 10.1038/nature14236.
- John F. Nash. Equilibrium points in  $n$ -person games. *Proceedings of the National Academy of Sciences*, 36(1):48–49, 1950. doi: 10.1073/pnas.36.1.48. URL <https://www.pnas.org/doi/abs/10.1073/pnas.36.1.48>.
- F. A. Oliehoek, M. T. J. Spaan, and N. Vlassis. Optimal and approximate  $q$ -value functions for decentralized POMDPs. *Journal of Artificial Intelligence Research*, 32:289–353, may 2008. doi: 10.1613/jair.2447. URL <https://doi.org/10.1613/jair.2447>.
- OpenAI, 2019. URL <https://openai.com/blog/openai-five-finals/>.
- R. Tyrrell Rockafellar. *Convex analysis*. Princeton Mathematical Series. Princeton University Press, Princeton, N. J., 1970.
- John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization, 2015a. URL <https://arxiv.org/abs/1502.05477>.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation, 2015b. URL <https://arxiv.org/abs/1506.02438>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.

- L. S. Shapley. Stochastic games\*. *Proceedings of the National Academy of Sciences*, 39(10):1095–1100, 1953. doi: 10.1073/pnas.39.10.1095. URL <https://www.pnas.org/doi/abs/10.1073/pnas.39.10.1095>.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Mach. Learn.*, 8(3–4):229–256, may 1992. ISSN 0885-6125. doi: 10.1007/BF00992696. URL <https://doi.org/10.1007/BF00992696>.
- Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising effectiveness of ppo in cooperative, multi-agent games, 2021. URL <https://arxiv.org/abs/2103.01955>.