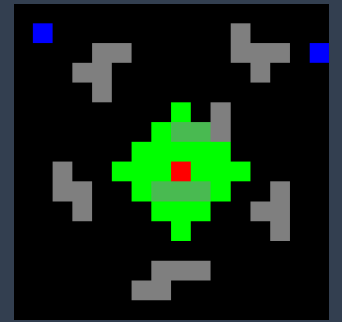




OÉ Gaillimh
NUI Galway



Exploring the Impact of Tunable Agents in Sequential Social Dilemmas

David O'Callaghan

ocallaghan1@gmail.com

Patrick Mannion

patrick.mannion@nuigalway.ie

School of Computer Science

National University of Ireland Galway

1. Abstract

2. Background

3. Related Work

4. Methodology

5. Experimentation

6. Results

7. Future Work

Abstract

- Reinforcement learning agents are most commonly trained to learn some fixed behaviour.
- However, if a different form of agent behaviour is required, the agent typically needs to be retrained.
- In this study, we use a method to design agents whose behaviours can be tuned after training using techniques from multi-objective reinforcement learning.
- We empirically show that the tunable agents framework enables easy adaption between cooperative and competitive behaviours in sequential social dilemmas.

Background

Multi-objective reinforcement learning

- Standard RL methods solve a sequential decision-making problem by optimising a single objective.
- Many real-world problems are multi-objective by nature and these objectives can be in conflict with each other.
- The MDP is extended to a multi-objective MDP, where the agent receives vectorized rewards.

$$\mathbf{V}^{\pi}(s) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \mathbf{r}_{t+k+1} \mid \pi, S_t = s \right]$$

Background

Multi-objective reinforcement learning continued...

- One approach to solving MOMDPs is to use a scalarisation function
- General form:

$$V_{\mathbf{w}}^{\pi}(s) = f(\mathbf{V}^{\pi}(s), \mathbf{w})$$

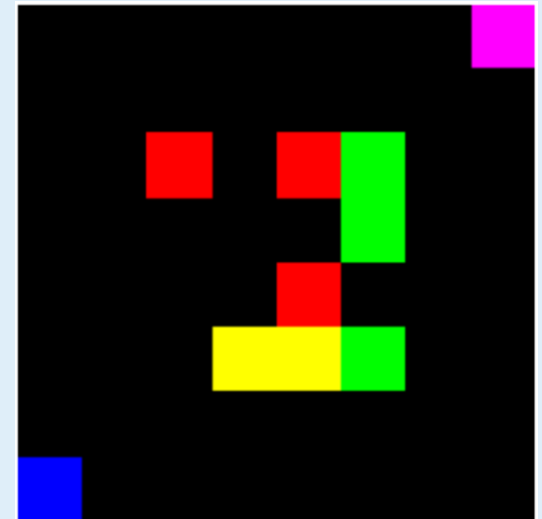
- Linear scalarisation:

$$V_{\mathbf{w}}^{\pi}(s) = \mathbf{w} \cdot \mathbf{V}^{\pi}(s)$$

- If \mathbf{w} is unknown in the learning stage, a set of policies can be learned using a range of values for \mathbf{w} .

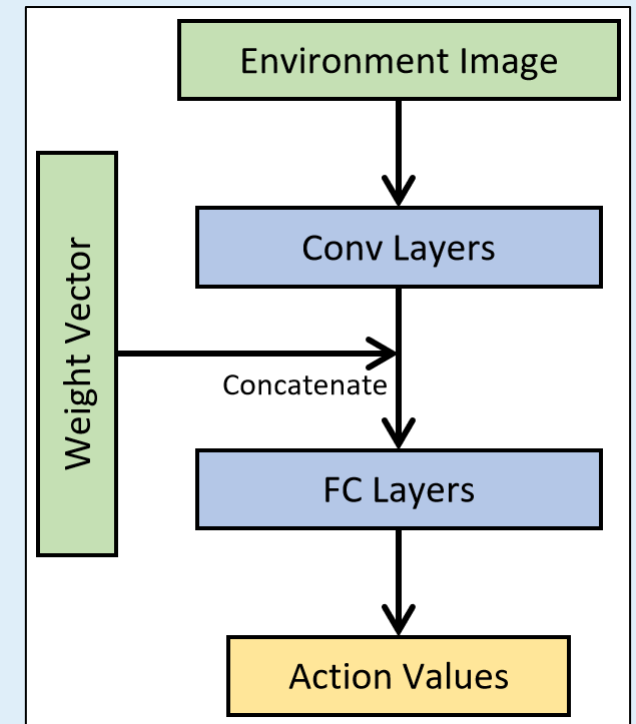
Related Work

- Källström and Heintz (2019) introduced an agent architecture for training agents that can prioritize a set of objectives in a multi-objective environment during runtime.
- The study showed that agents trained using this architecture could approximate policies of agents trained with fixed preferences over objectives.
- Used 2 gridworld environments with a single trainable agent in experiments.



Methodology

- The framework used for training tunable agents involves combining linear scalarization of rewards with the DQN algorithm.
- The neural network is conditioned on both the environment state and weight vector.
- The weight vector is re-sampled repeatedly during training.



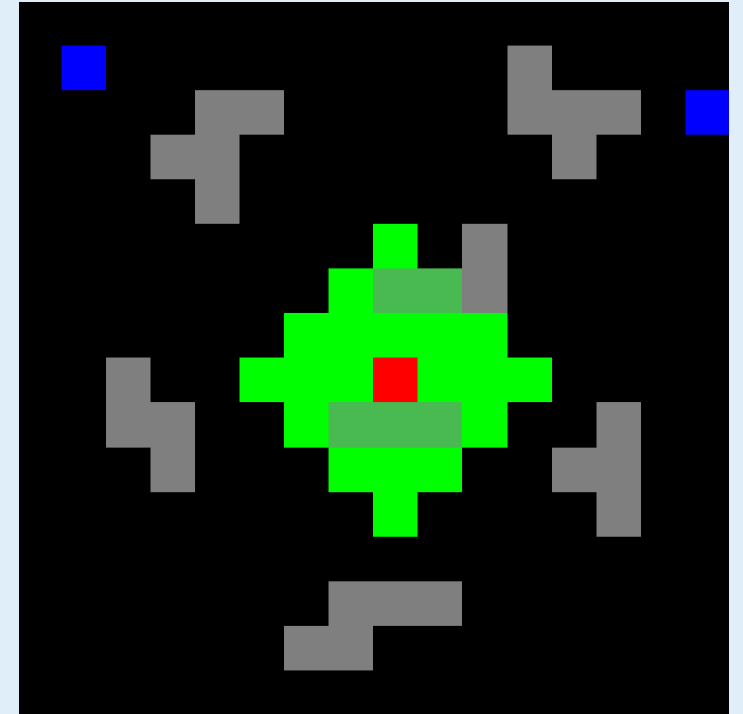
Methodology

The tunable DQN agent algorithm used in our work. *Adapted from Källström and Heintz (2019) for multi-agent settings.*

```
1 Initialise  $\epsilon$  to 1.0
2 Set  $T$  to the number of episodes to start training after
3 for  $episode \leftarrow 1$  to  $M$  do
4   Perform each command individually for agent  $i = 1, \dots, n$ 
   where  $n$  is the number of agents
5   Get the initial state  $s_i$  from the environment
6   Sample a set for preference weight vector ( $\mathbf{w}_i$ ) from the
   preference weight sample space
7   while  $s_i$  is not terminal do
8     Choose action  $a_i$  from state  $s$  using  $\epsilon$ -greedy policy:
      $a_i = \text{agent}[i].\text{get\_action}(s_i, \mathbf{w}_i, \epsilon)$ 
9     Take action  $a_i$  and observe reward vector  $\mathbf{r}_i$  and
     next state  $s'_i$ 
10    Store experience in replay memory:
      $\text{memory}[i].\text{append}(s_i, a_i, \mathbf{r}_i, s'_i, \mathbf{w}_i)$ 
11     $s_i \leftarrow s'_i$ 
12  end
13  Decay  $\epsilon$ 
14  if  $episode > T$  then
15     $\text{minibatch} \leftarrow \text{memory}[i].\text{sample}()$ 
16     $\text{agent}[i].\text{train}(\text{minibatch})$ 
17  end
18 end
```

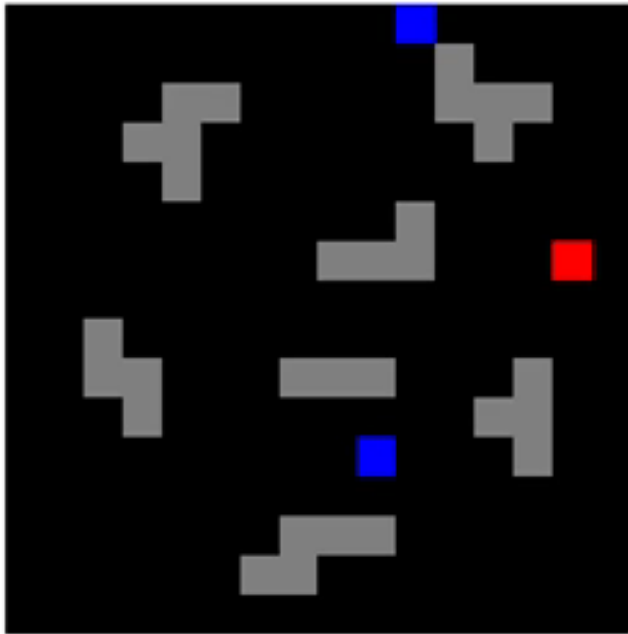

Experimentation

- In this study, we applied the tunable multi-agent framework to a multi-objective version of the Wolfpack environment from Leibo et al. (2017).
- 16x16 pixel-based gridworld. Multiple predator agents (blue) try to capture a prey (red).
- Team-capture if a predator is within the capture-radius (green area) when the other predator captures the prey. Lone-capture otherwise.
- Team-captures (cooperative) and lone-captures (competitive) are rewarded as separate objectives.

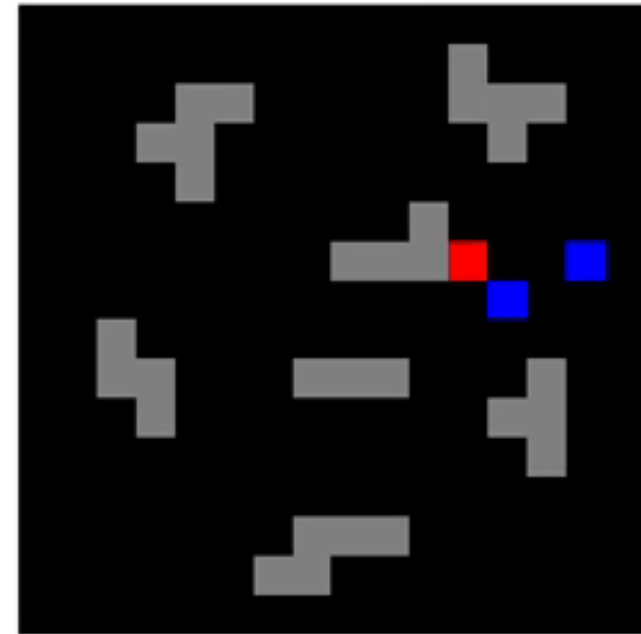


Results

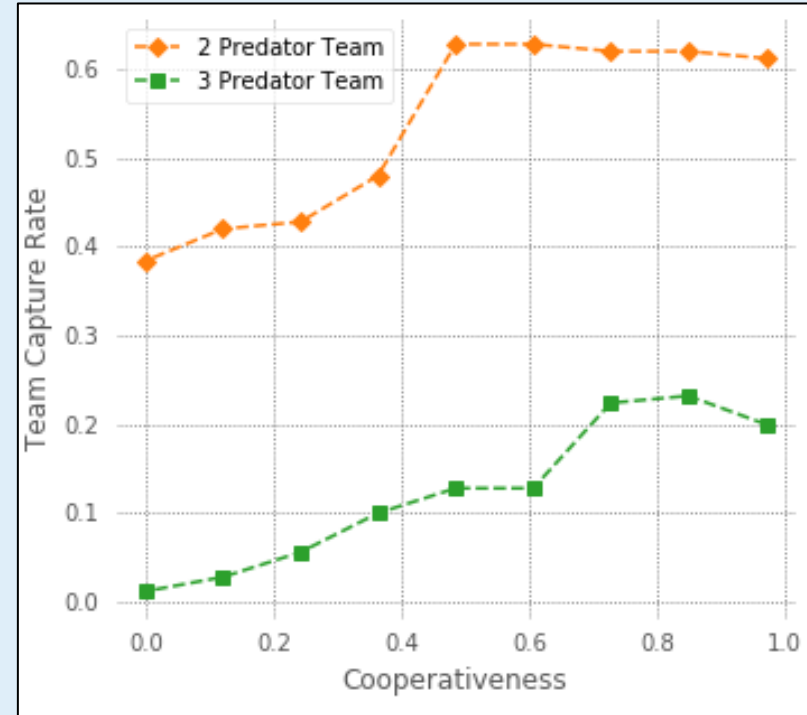
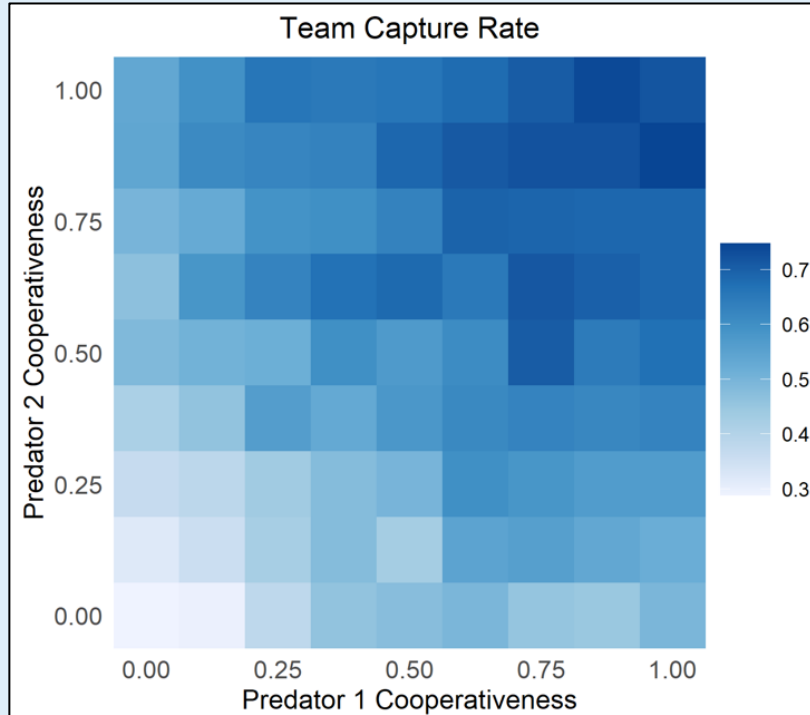
Competitive



Cooperative



Results



Results

Comparing tunable agents and fixed agents with empirical payoff matrices

		Pred. 2			
		C		D	
Pred. 1	C	0.792	0.792	0.592	0.296
	D	0.254	0.532	0.306	0.438

Tunable Agents

		Pred. 2			
		C		D	
Pred. 1	C	0.908	0.908	0.540	0.344
	D	0.340	0.600	0.436	0.468

Fixed Agents

Future Work

The tunable agent framework would be beneficial to any RL problem where there is some degree of uncertainty over the desired type of agent behaviour.

Wolfpack environment: One could train the prey to evade the predators and analyse the social dynamics.

Training scheme: One could easily use a different base RL algorithm to the DQN algorithm. Investigate the possibility of adapting the training methods for non-linear scalarisation.

Thank you for your attention.

Code available at:

<https://github.com/docallaghan/tunable-agents>