

CS8803 – Deep RL

FALL 2024

# Learning to Communicate in Multi-Agent Environments

---

**Presenters:** Mengying Lin, Nhi Nguyen, Wesley Ford, Nathan Chung (Group 12)

**Instructor:** Animesh Garg

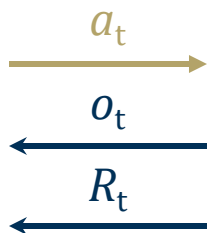


# Motivation

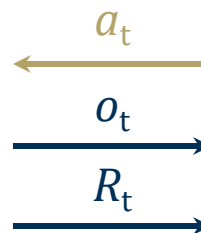
# Motivation and Main Problem



Agent

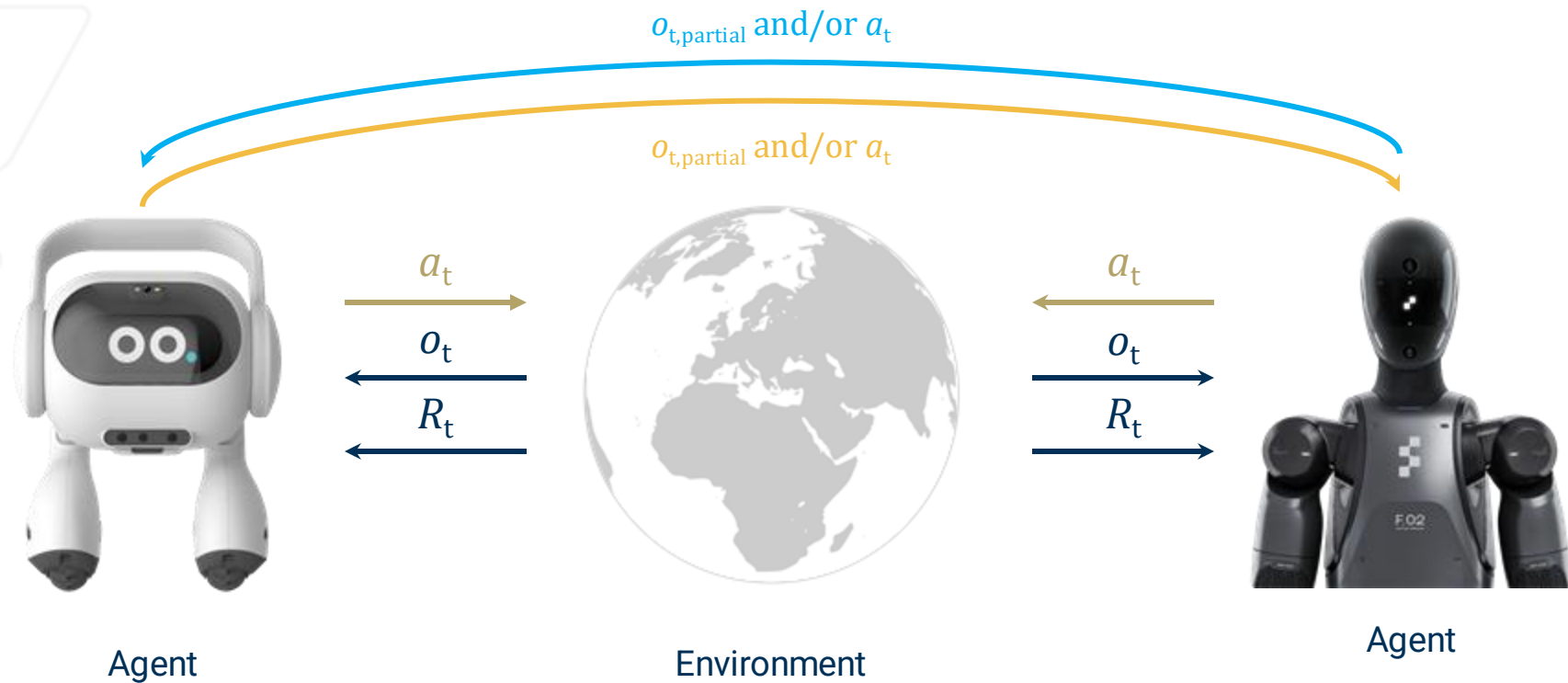


Environment



Agent

# Motivation and Main Problem

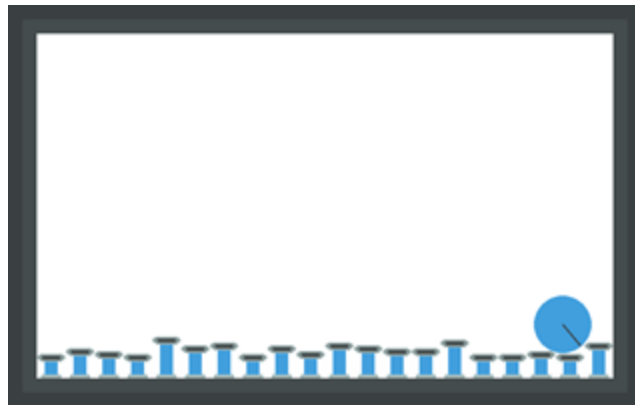




# Methodology

# Methodology

- We aim to investigate the impact of communication strategies and see how it impacts the performance overall. We decided to choose the **Pistonball environment** for its simplicity and abstract settings to focus on the coordinations performance.
- We chose **PPO** as the baseline to train multi-agents policy because it's a simple yet effective method. It also has been successfully used with competitive Multi-Agent Racing.



# Problem Settings

## Definitions and notations

Environment: Markov Decision Process (MDP)

$(S, A, P, r, p_0, \gamma)$

- A finite set of actions
- S finite set of states
- $P : S \times A \times S \rightarrow \mathbb{R}$
- $r : S \rightarrow \mathbb{R}$
- $p_0 : S \rightarrow \mathbb{R}$
- $\gamma \in (0, 1)$  (discount factor)

The Pistonball environment can be modeled as a Markov Decision Process (MDP).

Policy: target to maximize the value function  $V^\pi(s)$

$$V^\pi(s) = E_\pi[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_t = s]$$

Network: Actor-Critic

- Actor: output  $\pi_\theta(a|s)$  to choose next action.
- Critic: predict  $V^\pi(s)$

# Methodology

## Intuition

- Agents each use the observed state space to generate some emission that is used to augment the input to every agent's policies.
- Using a separate Neural Network, agents condense information from their observation into a single “communication” scalar. Each agent's communication scalar is shared with every other agent via a communication vector.
- Simple implementation and does not require much modification to the agents network and policy structure.

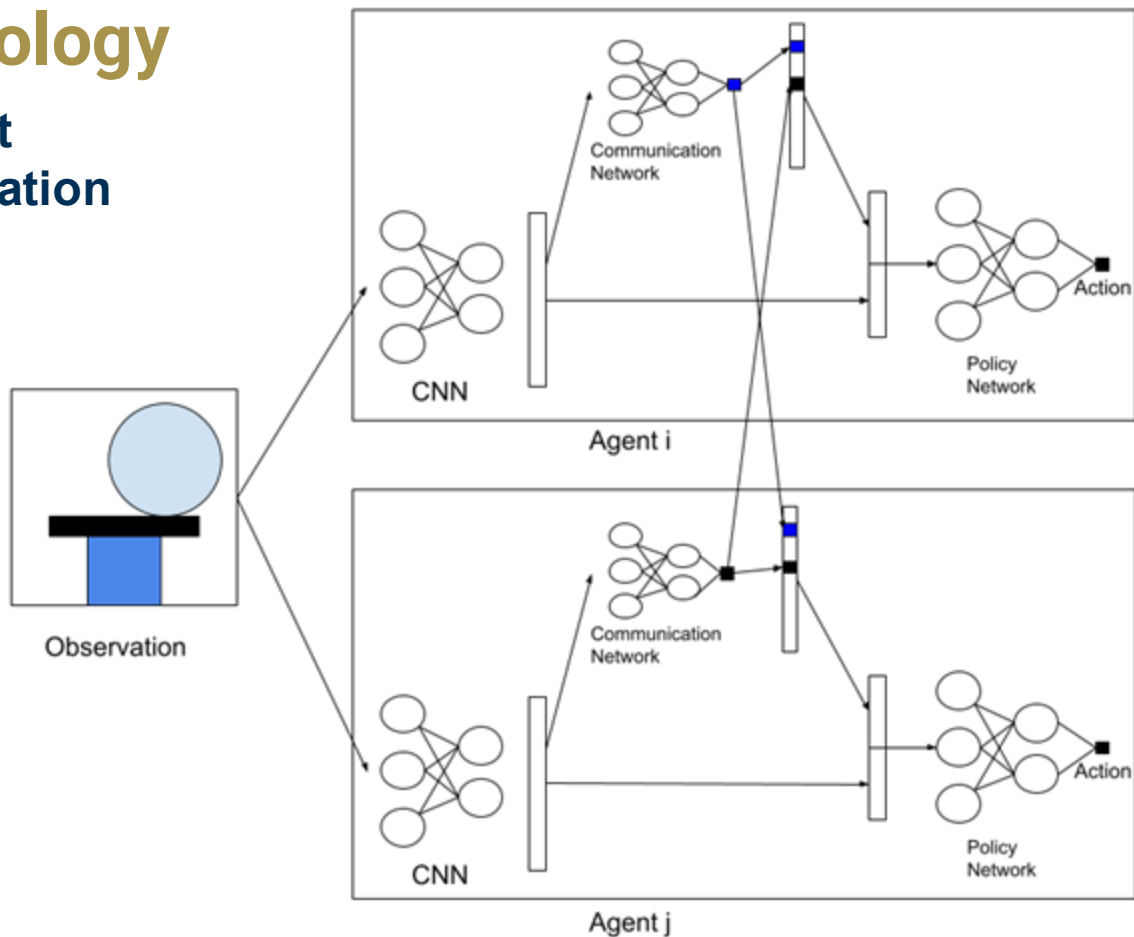


Multi-Agent Communication Vector



# Methodology

## Multi-Agent Communication Vector



# Methodology

## Intuition

- In multi-agents settings, agents must coordinates their actions to effectively complete the goal.
- Agents are receiving informations from multiple sources, it needs to cherry pick the suitable ones for the moment decision-making.
- As the number of agents increases, redundant information can lead to overload info and suboptimal performance.

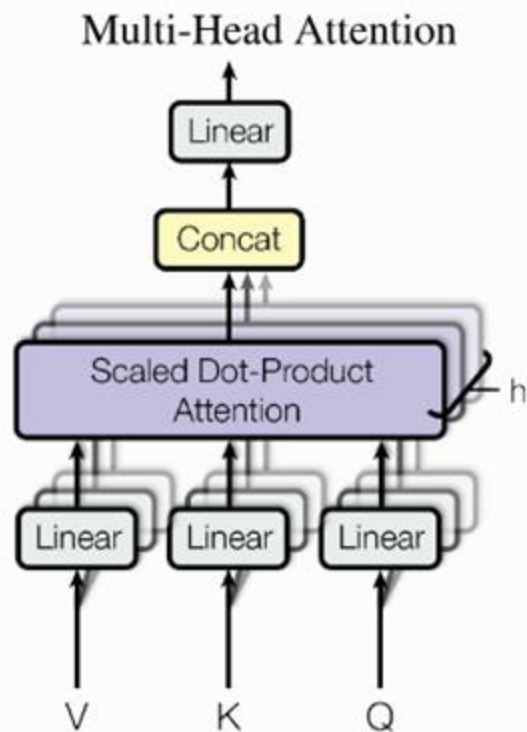


Multi-head attention for communication processing



# Methodology

## Multi-head attention for communication processing



$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

$$\text{Multihead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where  $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

# Methodology

---

**Algorithm 1** PPO with Multi-Agent Communication Vector

---

```
1: procedure PPO-MULTI-AGENT
2:   Initialize policy, value and communication networks  $\theta, \phi, \psi$  for each agent  $A_i$ .
3:   Initialize communication vector  $c_i$  for all agents.
4:   for each training iteration do
5:     for each environment step do
6:       for each agent  $A_i$  in  $\{A_1, \dots, A_N\}$  do
7:         Collect local observation  $s_i$ .
8:         Compute communication scalar  $c_i = \pi_\psi(s_i)$ 
9:         Compute communication vector by aggregating neighbors:
            
$$C_i \leftarrow \text{Aggregate}(\{c_j\}_{j \in \mathcal{N}(i)})$$

10:        Form augmented observation:  $\tilde{s}_i = [s_i, C_i, I_i]$ , conditioned on agent id:  $I_i$ 
11:        Select action  $a_i \sim \pi_\theta(a_i|\tilde{s}_i)$ .
12:      end for
13:      Execute joint action  $a = \{a_1, \dots, a_N\}$  and observe rewards.
14:    end for
15:    Update  $\theta$  and  $\phi$  using PPO.
16:  end for
17: end procedure
```

---

---

**Algorithm 2** PPO with Multi-Head Attention Communication

---

```
1: procedure PPO-MULTI-AGENT
2:   for each training iteration do
3:     for each environment step do
4:       for each agent  $A_i \in \{A_1, \dots, A_N\}$  do
5:         Collect local observation  $s_i$ ;
6:         Extract hidden features:  $h_i \leftarrow \text{Network}(s_i)$ 
7:         Concatenate hidden features with agent ID:
            
$$c_i \leftarrow \text{comm\_input}_i \leftarrow \text{Concat}(h_i, \text{id}_i)$$

8:         Compute communication vector using multi-head attention:
            
$$c_i \leftarrow \text{Attention}(\{c_j\}_{j \in \mathcal{N}(i)})$$

9:         Form augmented observation:  $\tilde{s}_i = [s_i, c_i]$ 
10:        Select action:  $a_i \sim \pi_\theta(a_i|\tilde{s}_i)$ 
11:      end for
12:      Execute joint action  $a = \{a_1, \dots, a_N\}$ ; observe rewards
13:    end for
14:    Optimize policy  $\theta$  and value function  $\phi$  using PPO
15:  end for
16: end procedure
```

---



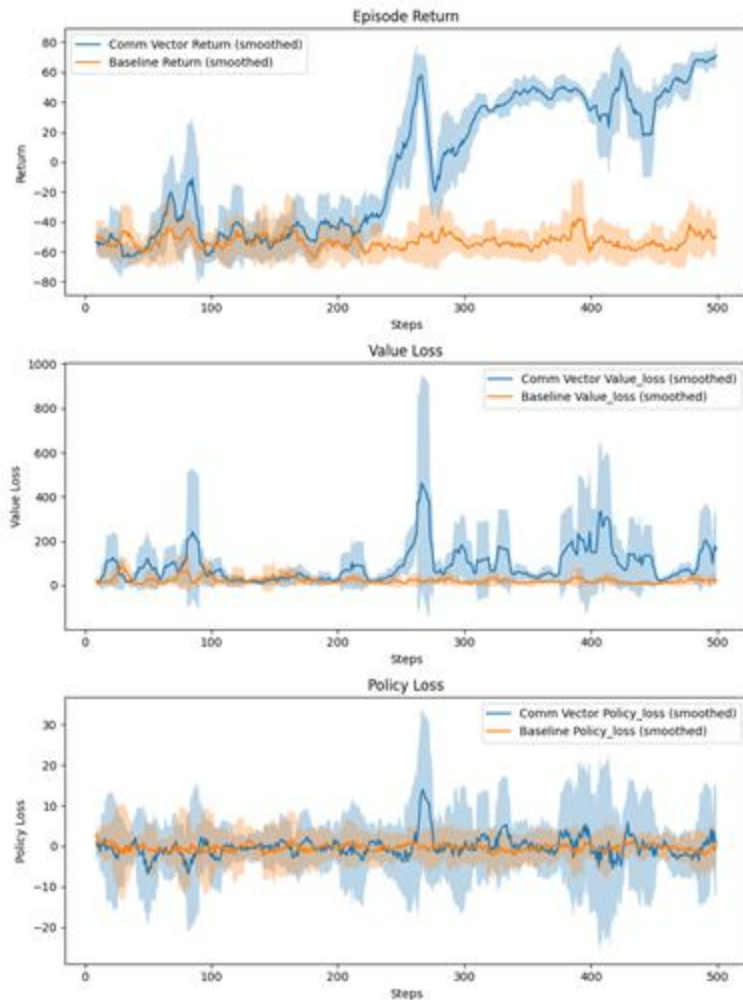
# Experimental Results

# Experimental Results

## Communication Vector

With a communication vector agent was able to converge to a solution within 500 steps.

High Variance Learning with communication. Learning what to communicate can create very large changes in the value function compared to a baseline.

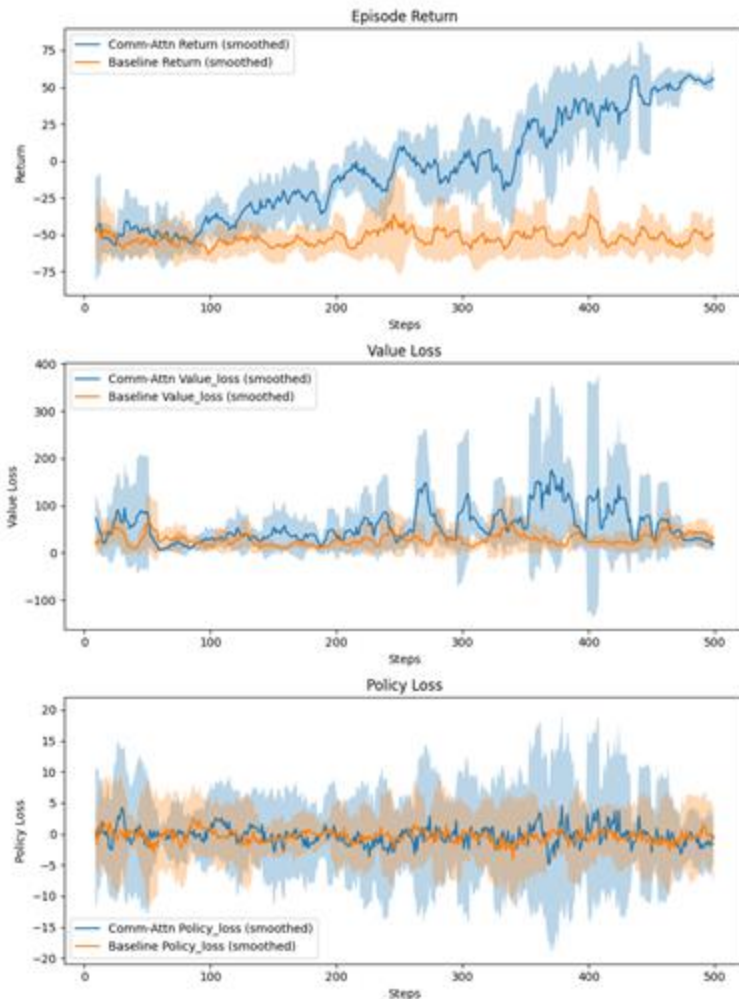


# Experimental Results

## Multi-head attention

Comm-Attn (blue) learn faster and performs much better than base-line

Since Comm-Attn learn from more complex input (communication vectors and attention), it learns slower from the start. It also has slightly higher variance in policy loss, showing that learning from communication is harder.

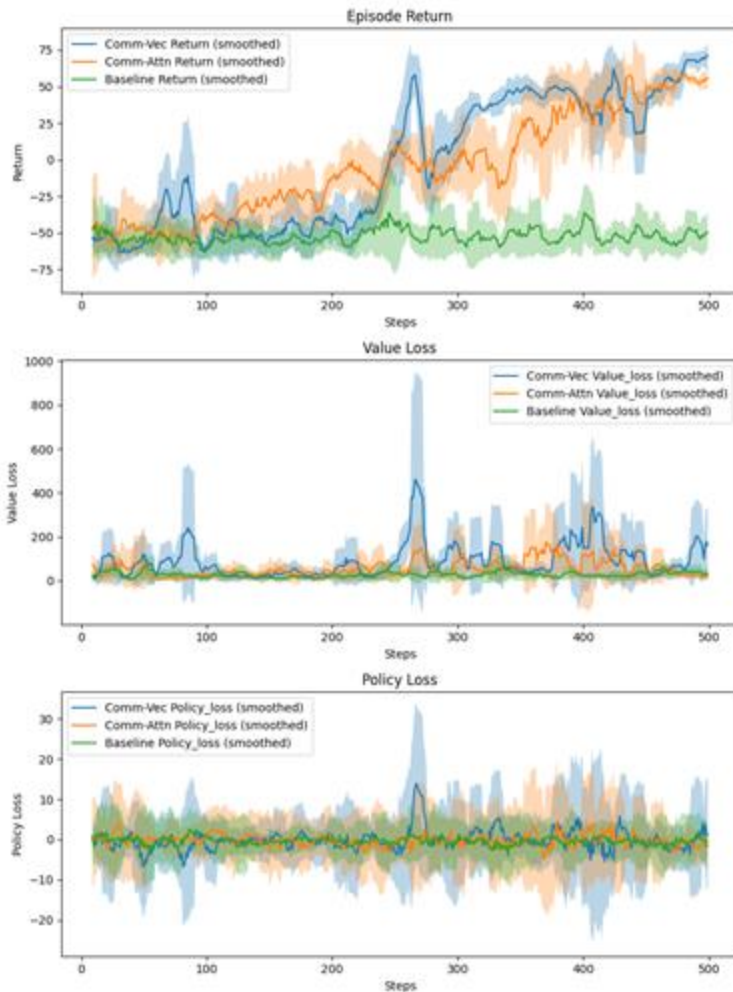


# Experimental Results

## Multi-head attention vs Communication Vector

Comm-Vec seems to prioritize aggressive learning, reflected in its higher variability in both value and policy losses but achieves better peak performance in returns

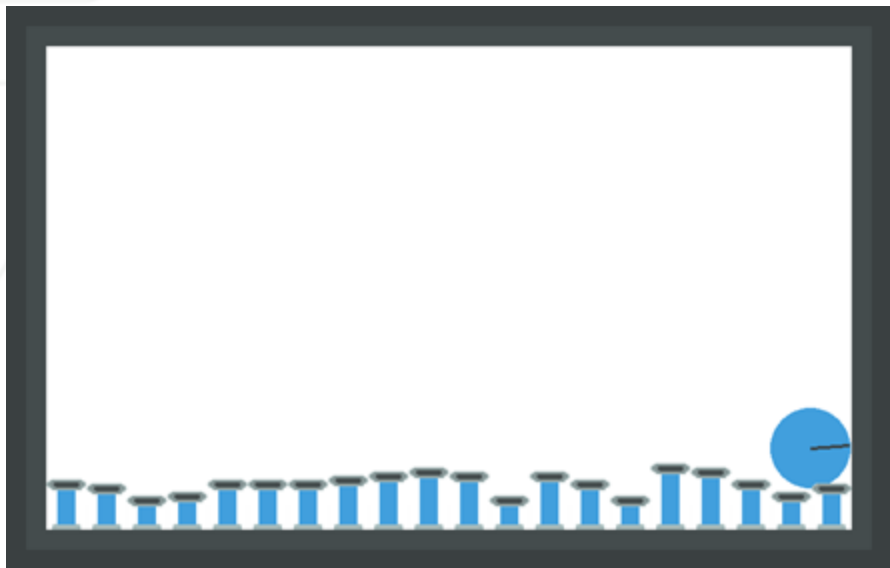
Comm-Attn trades off slightly lower returns for greater stability and consistency in the learning process, making it potentially more reliable.



# Experimental Results

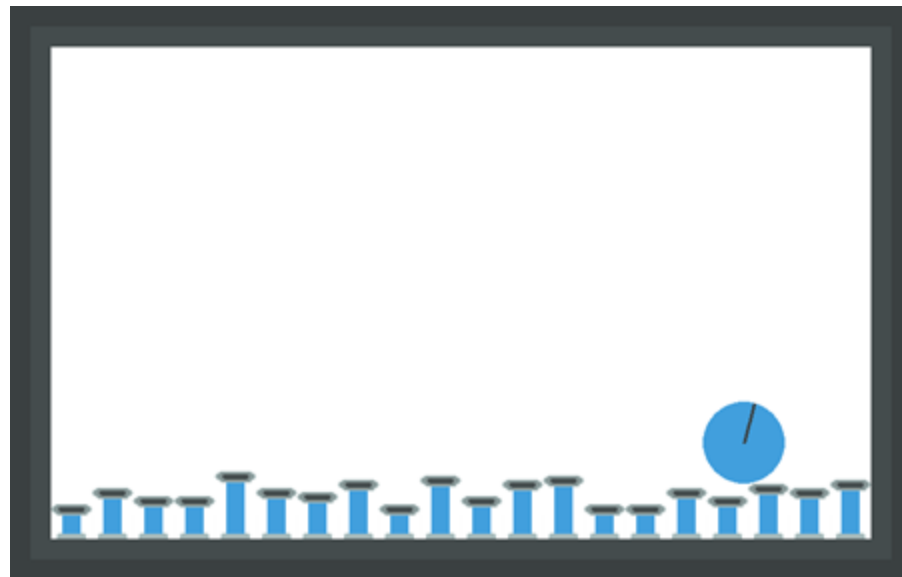
## Multi-head attention

### Our method



### Base-line

(Siqu Liu, et al 2019)



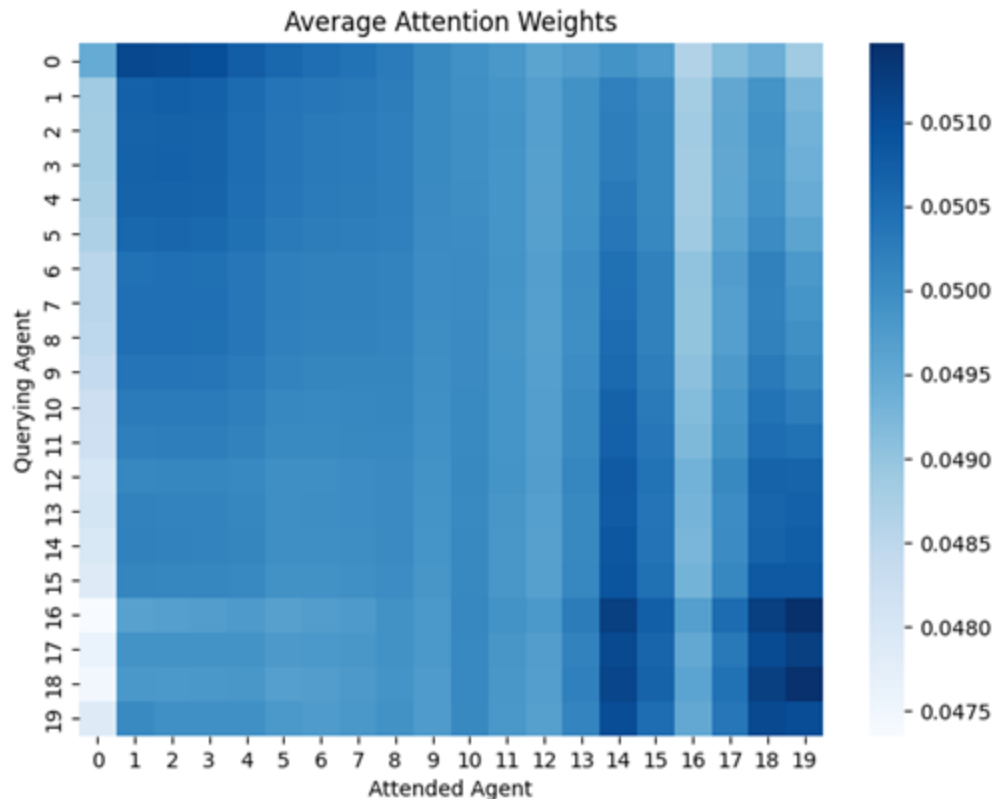
# Experimental Results

## Multi-head attention - Attention weights

Diagonal dominance: higher attention to themselves or immediate neighbors

Agents near the end 17-19 seems to give higher attention to specific attended agents

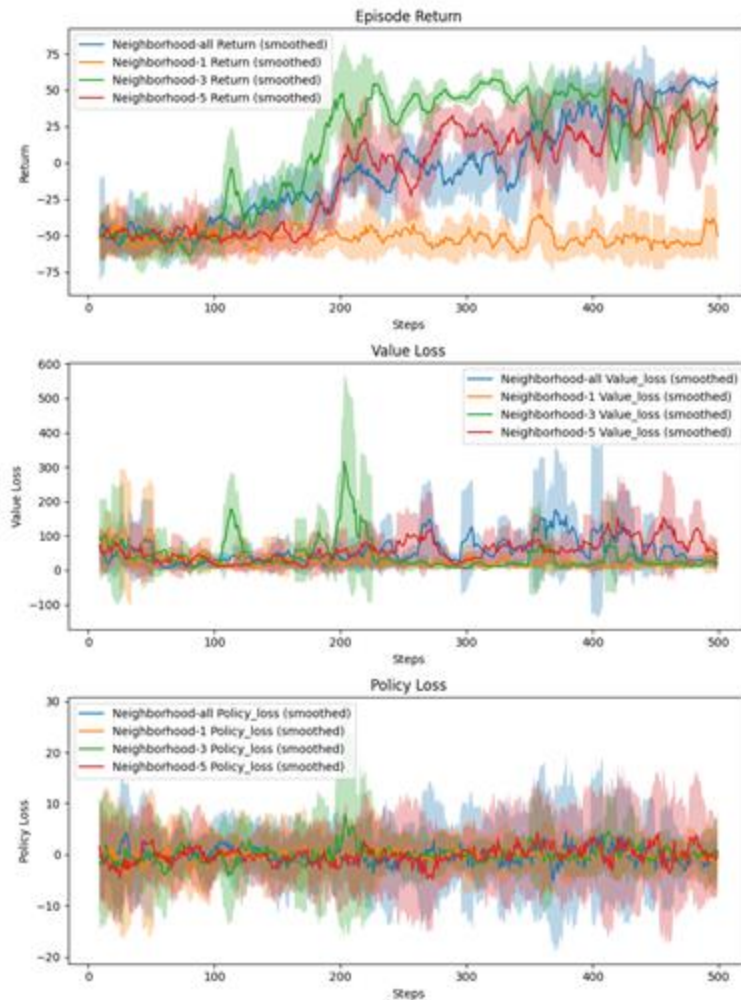
Low attention for far agents



# Experimental Results

## Multi-head attention - Neighbor range

Agents only need to focus on relevant, nearby neighbors (3-5) to perform well. Overloading them with information (all agents) or underloading them (1 agent) hurts performance.





# Conclusion

# Conclusion



Implementing communication can improve learning in collaborative multi-agent environments (which is observed in the results)



Simple methods such as using a communication vector and more complex methods using Multi Headed Attention were sufficient to solve simple environment



Future work would be to explore more complex multi-agent environments