# Predictive Auxiliary Learning for Belief-based Multi-Agent Systems

Qinwei Huang
*Electrical Engineering and Computer Sciences*
*Syracuse University*
Syracuse, USA
qhuang18@syr.edu

Stefan Wang
*Computer Science*
*University of Rochester*
Rochester, USA
swang170@u.rochester.edu

Simon Khan
*Air Force Research Laboratory*
simon.khan@us.af.mil

Garrett E. Katz
*Electrical Engineering and Computer Sciences*
*Syracuse University*
Syracuse, USA
gkatz01@syr.edu

Qinru Qiu
*Electrical Engineering and Computer Sciences*
*Syracuse University*
Syracuse, USA
qiqiu@syr.edu

*Abstract*—The performance of multi-agent reinforcement learning (MARL) in partially observable environments heavily relies on the effective aggregation of information from both observations and communications, as well as feedback signals like rewards. While most state-of-the-art multi-agent systems (MAS) primarily leverage rewards as the sole feedback for policy training, our research reveals that training agents to perform auxiliary predictive tasks can substantially enhance system performance. In this paper, we introduce BElief-based Predictive Auxiliary Learning (BEPAL), a novel approach that leverages auxiliary training objectives to provide additional information, aiding in policy learning. BEPAL follows the centralized training with decentralized execution (CTDE) paradigm. Each agent using BEPAL trains a belief model that predicts unobservable state information, such as other agents' rewards and motion directions, alongside its policy model. By enriching hidden states to represent critical environmental information that does not directly contribute to reward maximization, this concurrent auxiliary learning stabilizes the MARL process and improves its performance. We validate the effectiveness of BEPAL using two multi-agent cooperative games: the predator-prey game and the Google Research Football (GRF). In these scenarios, BEPAL outperforms existing methods, achieving an average 16% improvement in performance metrics and more stable convergence.

*Index Terms*—Reinforcement Learning, Multi-Agent System

## I. INTRODUCTION

Multi-Agent Systems (MAS) are extensively employed for tasks that require coordinated efforts among multiple autonomous agents, such as collaborative robotics and complex game environments. In MAS, the presence of multiple agents introduces greater uncertainty, increases computational demands for communication, and complicates the achievement of game-theoretic equilibria. To mitigate these challenges, Multi-Agent Reinforcement Learning (MARL) is utilized, enabling agents to learn to make optimal decisions through interaction with their environment and fellow agents.

Centralized Training with Decentralized Execution (CTDE) is a widely adopted strategy in MARL, designed to address the non-stationarity [1], [2] that arises when multiple agents learn and adapt simultaneously. CTDE uses a centralized value function during training, which is later decomposed into individual utility functions for decentralized execution, aligning agents' actions with both individual and collective objectives. Moreover, recent advancements [3]–[5] have incorporated communication systems into Centralized Critic with Decentralized Actor models to better manage challenges in partially observable environments. Previous works on communication-based Centralized Critic with Decentralized Actor [6]–[8] have demonstrated that integrating a communication system can significantly improve MARL performance. While these methods generally improve cooperation among agents, achieving high performance in complex tasks remains a challenge.

A DRL agent maps environment states to actions for maximum expected reward. The pursuit of enhanced state representation with richer information has recently become a key focus in addressing multi-agent challenges. For example, to obtain enriched state information, [9] integrates knowledge graphs with agent observations, [10] leverages attention mechanisms to enhance state by focusing on critical communication messages, and [11] combines predicted future events with current observation. However, these works adopt centralized control. They assume that the centralized controller has global observability and maps the observation plus auxiliary information (i.e., knowledge graph or prediction) directly to action. If their state enhancement method has trainable parameters, they are trained directly to maximize the rewards.

Many decentralized multi-agent games are partially observable systems. Agents can only observe their local environment and depend on communication to achieve global awareness. Additionally, these games are often non-Markovian from a single agent's perspective. The expected reward distribution depends not only on the target agent's current state but also on the states of other agents and hidden environment variables, which can be influenced by the target agent's past actions.

Existing approaches, such as [6], [8], [12], use Long short-term memory (LSTM) models to integrate an agent's historical observations and received messages into a hidden state. This hidden state needs to be comprehensive, capturing both current and past information about the environment. Decision-making based on this hidden state is considered to be Markovian. We refer to this hidden state as the *belief state* as it represents the agent's understanding of the environment. The belief state representation is trained through reinforcement learning using reward feedback.

Prior work has explored improving belief state representation by training agents to estimate the current locations of other agents, demonstrating enhanced policy performance [13]. However, this technique is not scalable to larger and more complex environment because it models the environments as a grid-based system and applies CNN to encode the observations. Moreover, to further enhance the belief representation, the agent needs to do more than estimating the location of other agents.

This work presents BElief-based Predictive Auxiliary Learning (BEPAL), a novel MARL technique that enhances agents' belief state representation by auxiliary predictive learning. BEPAL trains the agent to perform auxiliary tasks, such as predicting critical unobservable or partially observable information like rewards and the locations of other agents, to improve the performance of MARL. BEPAL supplements each agent with a belief decoder and a set of auxiliary objective functions, which train the decoder to transform the agent's hidden state into a neuro-symbolic representation of its understanding of the global environment. This representation may contain predictive information of other agents' locations, moving directions, and future rewards. While they may not directly affect the reward of the target agent, they influence the outcome of collaboration activities, and, consequently, the overall system reward. Our experimental results demonstrate that enriching the belief state representation by training agents to perform these auxiliary predictive tasks significantly improves performance in multi-agent distributed systems where agents make local decisions.

Additionally, we employ a multi-head graph attention neural network (GAT) [14] as the observation encoder to handle complex observations in non-grid based environment. The complexity of this encoder is proportional to the size of the agent's observations and independent to the environment's size. Unlike convolutional neural networks (CNNs), which operate on Euclidean vector space, the GAT based encoder accommodates arbitrary precision in object locations and non-Euclidean inputs. With these enhancements, BEPAL not only surpasses previous models in agent performance but is also adaptable to new environments that were unsupported by some of the earlier models.

## II. RELATED WORKS

In multi-agent settings, each agent can perceive other agents as part of their state observations, resulting in a dynamic and non-stationary environment [15], [16]. Centralized Training

with Decentralized Execution (CTDE) effectively addresses these challenges by maintaining strong generalization abilities and resolving non-stationarity better than centralized control methods like QMIX [17]. A key advancement in this area is the use of a centralized critic in actor-critic algorithms, which facilitates decentralized policy learning, as demonstrated by Multi-Agent Deep Deterministic Policy Gradient (MADDPG) [18]. Despite its strengths, MADDPG struggles with inadequate experience collection, impeding performance and complicating convergence. COMA [19] further improves centralized critic with counterfactual baselines to enhance credit assignment and optimize decentralized policies in cooperative multi-agent systems.

While these approaches mark critical improvements in MARL, the introduction of communication systems among agents has proven essential for further enhancing performance. However, the complexity of applying joint value functions increases with the number of agents, making centralized training more challenging [20]–[23]. Innovations like Reinforced Inter-Agent Learning (RIAL) and Differentiable Inter-Agent Learning (DIAL) [24] introduced message-passing frameworks based on local observations and actions. CommNet [12] integrated a centralized communication channel, while IC3Net [6] added a communication gate to optimize timing. However, these models often overlook the content of messages, relying on simple averaging and equal weighting, which may not capture the full dynamics of the environment. TarMAC [8] and MAGIC [7] address this by employing attention mechanisms to selectively weight and direct message passing, improving the relevance of shared information. PR2 [25] approximates opponents' conditional policies using variational Bayes methods to iteratively update the agent's own strategies.

Despite these advancements, many MARL approaches overlook key aspects of team dynamics and environmental context, such as the varying positions and rewards of other agents. In DRL, auxiliary objectives have proven effective for enhancing state representation by predicting rewards or future states [26]–[29]. In MARL, however, it is essential to integrate information about other agents' states to improve the overall understanding of the environment. BAMS [13] uses CNNs on grid maps to provide additional feedback, though it is limited to grid-based environments and incurs additional computational expense.

Our approach differs by extracting key information within MAS for auxiliary prediction objectives to enrich the agents' belief states. By integrating predictions of rewards from other agents, we enable each agent to build a more accurate understanding of the overall multi-agent system, leading to improved cooperation and performance.

## III. PRELIMINARIES

### A. Dec-Partially Observable Markov Decision Processes

We investigate a cooperative multi-agent game within the framework of the Decentralized Partially Observable Markov Decision Process (Dec-POMDP), defined as the tuple $\langle N, S, P, \mathcal{R}, \mathcal{O}, \mathcal{A}, Z, \gamma \rangle$. Here, $N$ denotes the number of agents; $S$ is the finite state space; $P(s' \mid s, a) : S \times$

$\mathcal{A} \times S \rightarrow [0,1]$ represents the state transition probabilities; $\mathcal{A} = [\mathbf{A}_1 \ldots \mathbf{A}_N]$ constitutes the finite set of actions, with $\mathbf{A}_i$ indicating the local actions $\mathbf{a}_i$ available to agent $i$; $\mathcal{O} = [\mathbf{O}_1 \ldots \mathbf{O}_N]$ encapsulates the finite set of observations governed by the observation function $Z : S \times \mathcal{A} \rightarrow \mathcal{O}$; $\mathcal{R} : S \times \mathcal{A} \rightarrow \mathbb{R}^N$ is the reward function; and the constant $\gamma \in (0, 1]$ is a discount factor. In each time step $t$, agent $i$ selects action $\mathbf{a}_i^t$, and receives reward $r_i^t$ and observation $o_i^t$. Agent $i$ aims to maximize its discounted reward $R_i = \sum_{t=0}^{T} \gamma^t r^t$.

### B. Multi-head Graph Attention Networks

Graph Attention Networks (GATs) leverage a self-attention mechanism to refine node representations within a graph. Each node $i$ evaluates the importance of its neighboring nodes $j$, where $j \in \mathcal{N}_i$ and $\mathcal{N}_i$ is the set of one-hop neighbors of $i$. The attention coefficient $\delta_{ij}$ is calculated as follows:

$$w_{ij} = LeakyReLU(a^T[Wf_i \| Wf_j]) \quad (1)$$

$$\delta_{ij} = \frac{e^{w_{ij}}}{\sum_{k \in \mathcal{N}_i} e^{w_{ik}}} \quad for \; j = 1, 2, \ldots, |\mathcal{N}_i| \quad (2)$$

where $[\cdot \| \cdot]$ denotes concatenation, $f_i$ and $f_j$ are features of nodes $i$ and $j$, $W$ is a trainable weight matrix, and $a$ is a trainable weight vector.

Using the attention coefficient, node $i$ updates its feature vector by aggregating information from its neighbors:

$$\vec{f}_i = LeakyReLU \left( \sum_{j \in \mathcal{N}_i \bigcup \{i\}} \delta_{ij} W f_j \right) \quad (3)$$

To encapsulate more complex relationships among nodes, a multi-head GAT configuration is employed. Each attention head computes its own attention coefficient $\delta_{ij}^k$ and produces an updated feature vector $\vec{f}_i^k$, where $k$ denotes the head index.

The GAT layer updates a node's representation using the features of its one-hop neighbors. To incorporate the features from N-hop neighbors, an N-layer of GAT is required.

### IV. METHOD

Our target system comprises $N$ agents, each equipped with broadcasting-based communication capabilities. From each agent's perspective the environment is modeled as a POMDP. At each time step $t$, an agent $i$ receives a local observation $o_i^t$ and inter-agent messages $m_{ij}^{t-1}$ from other agents $1 \leq j \leq N$. Using this information, the agent updates its hidden state $h_i^t$. A policy $\pi_i$ then maps the hidden state $h_i^t$ to a probability distribution over the action space. Each action consists of two components: a communication action that determines whether message broadcasting is gated or allowed, and a movement action that dictates how the agent plays the game. The actual action $a_i^t$ is sampled according to $a_i^t \sim \pi_i$.

In this work, we focus on homogeneous agents and adopt a model sharing approach, where the same model is used by different agents with their local observations $o_i$. The model enables the agents to maintain their local hidden states $h_i$ and select local actions $a_i$. A shared model typically can be trained

faster than the agent specific model because $N$ times more training data can be collected from each game.

Each agent also uses the hidden state $h_i^t$ to update its beliefs $b_i^t$ and $p_i^t$ about other agents' motion and rewards, respectively. Maximizing the accuracy of these beliefs serves a set of auxiliary training objectives. The subsequent subsections will detail the operational flow of the framework, describe the feature extraction models for the received input, elaborate on the functionalities of the belief decoder module, and delineate the loss functions used to train the system.

The architecture of the BEPAL model is depicted in Fig 1. It consists of a feature extractor, an LSTM, a policy network and a belief decoder network.

### A. Feature Extractors

For agent $i$ at timestamp $t$, the inputs, which consisting of local observations $o_i^t$ and inter-agent messages $m_{ij}^{t-1}$, are processed by the feature extractor. The observation $o_i^t$ is represented as a graph $G(Vertex, Edge)$, where each vertex represents an object observed by the agent except vertex 0, which presents the agent itself. We consider only the relationship between the observed objects and the agent, hence the graph has a star connection, $Edge = \{(edge_0, edge_i) | \forall j, 0 < j \leq N\}$. Each vertex in the graph is associated with a feature vector.

The feature extractor processes the agent's observations using a two-layer Graph Attention Network (GAT), trained to assign higher attention weights to objects containing critical information, such as goals or enemies, to enable more informative feature extraction. The first layer employs 3 attention heads while the second layer has a single attention head. We select only the GAT output for node 0, $\vec{f}_0$, as it represents the feature of agent $i$ updated with aggregated information from its observations. $\vec{f}_0$ is further processed by a fully connected layer, and the output, $e_i^t$, serves as the encoded observation of agent $i$.

In addition to local observations, communication messages exchanged between agents play a crucial role in decision-making. These messages encapsulate historical information gathered by each agent and are shared with their teammates. This approach is inspired by frameworks utilized by TarMAC and IC3Net, where the hidden state held by local LSTM also serves as the communication message. Not all received messages are equally important to an agent's decision-making, highlighting the importance of efficient information extraction. To address this, an attention mechanism, similar to the GAT used for observations, is employed for effective message aggregation. The query, key and value ($q_i^t$, $k_j^t$, $u_j^t$) of the attention model are generated from the agent's hidden state and received messages by passing through a Linear layer respectively, $q_i^t = W_q h_i^t$, $k_j^t = W_k m_j^{t-1}$, $u_j^t = W_u m_j^{t-1}$, where $W_q$, $W_k$ and $W_u$ are trainable weight matrices. The aggregated message $c_i^t$ is then computed following Equations
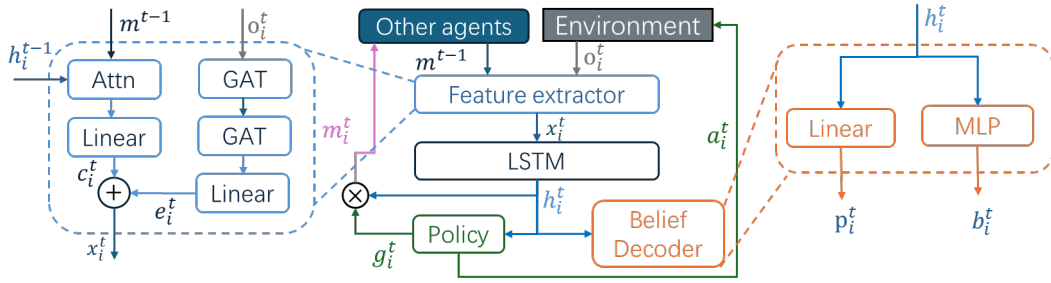
Fig. 1: The BEPAL-MAS architecture.

4 and 5:

$$\alpha_i^t = \text{SoftMax} \left[ \frac{\left( q_i^{t^T} k_1^t \right)}{\sqrt{(d_k)}}, \dots, \frac{\left( q_i^{t^T} k_j^t \right)}{\sqrt{(d_k)}}, \dots, \frac{\left( q_i^{t^T} k_N^t \right)}{\sqrt{(d_k)}} \right] \tag{4}$$

$$c_i^t = \sum_{j=1}^N \alpha_{ij}^t u_j^t \tag{5}$$

### B. Hidden State Update and Decision Making

The aggregated observation ($e_i^t$) and message ($c_i^t$) are summed up as $x_i^t$ and fed into an LSTM network. This LSTM enables the agent to retain historical information, making it easier to learn long-term goals and facilitating better policy learning by approximating the hidden state of the latent Markovian Process of the environment. The agent's hidden state $h_i^t$ and cell state $s_i^{t+1}$ are updated using Equation 6.

$$h_i^{t+1}, s_i^{t+1} = lstm(e_i^t + c_i^t, h_i^t, s_i^t) \tag{6}$$

Based on the hidden state, the agent chooses actions using a trained policy network within the actor-critic framework. The actor network, denoted as $actor()$, consists of a fully connected linear layer with LogSoftmax activation function. It maps the hidden state to a vector of action probabilities, $\pi(a|h)$, also known as policy. The action taken by agent $i$ at timestamp $t$, denoted as $a_i^t$, is sampled from the policy distribution $a_i \sim \pi(a_i^t|h_i^t)$.

At each step, an agent must make two decisions, (1) to select a movement action to advance the game objectives, and (2) to decide whether to broadcast its hidden state. These actions are denoted as $a_i^t$ and $g_i^t$ respectively. The corresponding policies are denoted as $\pi_a$ and $\pi_g$. The outgoing message $m_i^t$ for agent $i$ at timestamp $t$ is computed as the element-wise product of hidden state and gate action: $m_i^t = h_i^t \odot g_i^t$.

The critic model, denoted as $critic()$, is a single-layer fully connected network. It gives agent $i$'s local estimation, $v_i^t$, of the total discounted future rewards.

$$a_i^t, g_i^t \sim actor(h_i^t), \quad v_i^t = critic(h_i^t) \tag{7}$$

### C. Auxiliary Predictive Tasks

In addition to the policy network, the hidden state $h_i^t$ is also used as the input for a belief decoder, which performs two auxiliary tasks: reward prediction and motion prediction.

The reward prediction is handled by a single layer fully-connected linear network $p_i^t = P_\theta(h_i^t)$. The output, $p_i^t \in R^N$, is a vector of expected rewards. The $j$th entry in $p_i^t$ gives agent $i$'s estimation of the discounted future reward of agent $j$. The motion prediction uses a 2-layer fully connected linear network, $b_i^t = B_\theta(h_i^t)$. The output, $b_i \in R^{N \times M}$, predicts $M$ motion features of the $N$ agents. The specific motion features to be predicted depend on applications. The same decoder is used by all agents to decode their local states. Again a shared model will maximize the training effectiveness.

We hypothesize that the more effectively an agent understands its environment, the better its decisions will be. A thorough understanding of the environment requires efficient feature extractions from observations and messages, as well as accurate state updates that retain useful information and discard irrelevant history. Relying solely on deep reinforcement learning (DRL) to train the feature extractor and LSTM is insufficient, as the game's reward signal does not provide direct feedback on state representation. We propose that by training the agent to accurately predict the future motion and rewards of its teammates from its hidden state $h_i^t$, we can improve the feature extractor and LSTM, leading to a more refined representation of the hidden state.

Another advantage of these auxiliary prediction tasks is that they encourage the retention of teammate information within the agent's hidden state. Since this information may not directly influence the agent's reward, it could be overlooked if the system is trained solely using DRL. However, retaining such information is crucial for making collaborative decisions, such as communication. The auxiliary predictive tasks are only performed during the training. The belief decoder is not required during the testing.

### D. Loss Functions

We employ centralized training and distributed execution, with all agents sharing the same BEPAL model and being trained together.

The training loss for each agent is composed of two main components: the actor-critic loss $L_{RL}$ and the auxiliary objectives losses $L_{aux}$, represented as $Loss = L_{RL} + \lambda L_{aux}$, where $\lambda$ is a hyper-parameter. The auxiliary loss is derived by comparing the agents' beliefs with the ground truth, specifically $\overline{b^t}$

and $\overline{p^t}$. The Mean Squared Error (MSE) is utilized to quantify the auxiliary objectives loss, as outlined in the equation:

$$L_{aux} = \sum_t^T \sum_i^N \left( \mu MSE(\overline{b^t} - b_i^t) + \nu MSE(\overline{p^t} - p_i^t) \right) \quad (8)$$

where $N$ represents the number of agents, and $T$ denotes the total number of time steps. $\mu$ and $\nu$ are additional hyperparameters indicating which belief type (motion or reward) is more important. During training, a central controller is needed to monitor all agent motion and observations and generates the ground truths.

The actor-critic loss is $L_{RL} = L_{actor} + \beta L_{critic}$, where $\beta$ is an importance hyper parameter. The critic loss aims to minimize the errors between the estimated state value and the discounted return:

$$L_{critic} = \sum_t ||r_i^t + \gamma V(h_i^{t+1}) - V(h_i^t)||^2 \quad (9)$$

where $r_i^t = \mathcal{R}(h_i^t, a_i^t)$ specifies the reward for agent $i$ at time $t$, and $V()$ denotes the value estimation by the critic model. The actor loss aims to maximize the expected cumulative reward by updating the policy parameters in the guided direction of critic. The actor network is updated via policy gradient as:

$$\nabla_\theta J(\theta) = \sum_t \nabla_\theta \log \pi_\theta(a_i^t|h_i^t)[r_i^t + \gamma V(h_i^{t+1}) - V(h_i^t)] \quad (10)$$

where $\pi_\theta$, parameterized by $\theta$, represents the policy function. The BEPAL-MAS model is updated with the average gradient of all agents, aiming to achieve consistent improvements across the system.

## V. EXPERIMENTS

### A. Evaluation Environments

To evaluate the performance of the BEPAL-based multi-agent system (BEPAL-MAs), we conducted experiments in two distinct environments: Predator-Prey and the complex Google Research Football (GRF) environment.

**Predator Prey (PP)** is a widely used MARL benchmark where $N$ predators (agents) with a limited vision range ($3 \times 3$) search for a stationary prey in a grid-based environment. Both the agents and prey are randomly initialized in the environment. The agents receive a -0.05 penalty until they catch the prey, after which they no longer receive a penalty. The game concludes when all predators have reached the prey and the goal is to complete the game with the fewest steps. We tested environments with three complexity levels by varying the map size and the number of randomly placed obstacles that agents cannot pass through as shown in Fig. 2.

**Google Research Football (GRF)** is a physics-based 3D soccer simulator designed for reinforcement learning. GRF offers 19 actions (e.g., movement, kicking, dribbling). The field dimensions are $2 \times 0.82$. We evaluated the models in the academy scenario 3 vs. 1 setting, where three attackers controlled by our model face off against one defender and one goalkeeper controlled by the built-in AI. To simulate a partially observable environment, we limited the agents' vision range

TABLE I: Average steps comparison in predator prey environment without obstacle

|  | N=5 m=12 Max Steps = 40 | N=7, m=16 Max Steps = 60 | N=5, m=20 Max Steps = 80 |
|---|---|---|---|
| IC3Net | $32.90 \pm 5.22$ | $53.54 \pm 7.51$ | $71.34 \pm 8.16$ |
| TarMAC | $34.59 \pm 10.48$ | $44.68 \pm 9.19$ | $79.94 \pm 5.98$ |
| MAGIC | $26.89 \pm 9.73$ | $46.04 \pm 17.94$ | $78.36 \pm 9.49$ |
| BAMS | $21.64 \pm 4.26$ | $30.76 \pm 6.98$ | $53.79 \pm 16.49$ |
| Ours | $\mathbf{18.23 \pm 2.41}$ | $\mathbf{29.38 \pm 5.42}$ | $\mathbf{43.77 \pm 8.48}$ |

to a radius of 0.5. The game concludes when a goal is scored, the ball goes out of bounds, or the possession changes. Each agent receives a reward of +1 for scoring a goal, 0 otherwise. Performance was measured by the average success rate, with a higher rate indicating superior model performance.
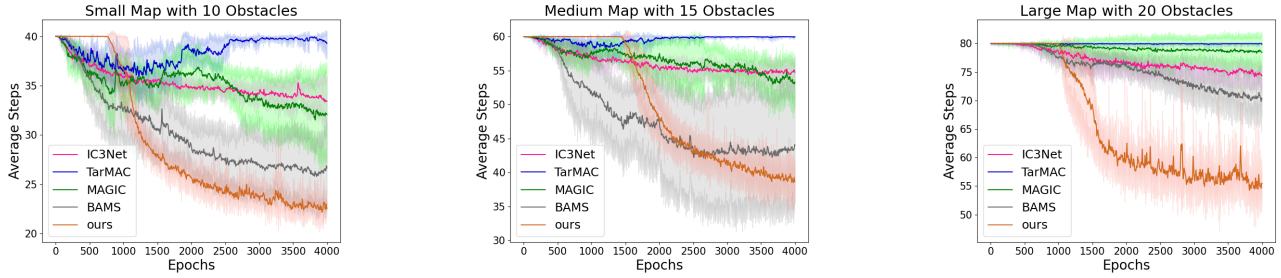
### B. Experiment setup

We optimized our network using RMSprop [30], with a learning rate of 0.001 and a smoothing constant of 0.97. Entropy regularization was applied with a coefficient of 0.01. Coefficient of critic loss $\beta$ is set to 0.05 and coefficient of auxiliary objective loss $\lambda = (\mu, \nu)$ is set to $(0.05/N, 0.5)$. The discounting factor for both games are set to be 1 since the outcomes of both games are given at the end. These hyperparameters were selected by hand and found to work well in practice, although a more systematic grid search could be used in future work. Code will be released upon publication.

We compare BEPAL with four baseline MARL models that leverages inter-agent communications with individual controllable communication gate: IC3Net, which employs message gating; TarMAC, which uses both message gating and attention; MAGIC, which leverages GNN-based models for message generation and passing in a centralized execution setup; and BAMS, which uses a grid-shaped belief map as one feedback channel. All experimental results are averaged over 5 seeds.

### C. Main Results for Predator Prey Game

*1) Performance Comparison:* Table I compares the BEPAL with four baseline algorithms in terms of average number of steps required to complete a game in an environment without obstacles. In the small (i.e., $m = 12$) and medium maps (i.e., $m = 16$), BEPAL-MAS performs comparably to the best baseline in terms of average steps, however, with a lower standard deviation. On the larger map (i.e., $m = 20$), it achieves a 18.62% reduction in average steps and a 50% reduction in the standard deviation compared to the best baseline.

We then introduced obstacles to create a more complex environment. Figure 2 shows the average number of steps required for all 5 models to complete the game. Compared to the baselines, our method demonstrated superior performance, reducing average step counts by 15.42%, 10.82% and 21.67% in the small, medium and large maps, compared to the best baseline model, respectively (Figure 2a-2c). Moreover, our method achieved more stable learning, with a 60.91% and

(a) Small map, 5 agents, 10 obstacles  (b) Medium map, 7 agents, 15 obstacles  (c) Large map, 5 agents, 20 obstacles

Fig. 2: Performance comparison for the Predator-Prey games with obstacles is conducted across maps of different sizes: Small ($12 \times 12$), medium ($16 \times 16$) and large ($20 \times 20$)

TABLE II: Ablation study on BEPAL

| Environment | GRF | PP with Obstacle |
|---|---|---|
| Evaluation Metric | Success Rate | Average Steps |
| BEPAL | 83.64% | 22.21 |
| BEPAL w/o Reward Prediction | 72.97% | 24.52 |
| BEPAL w/o Motion Prediction | 19.14% | 36.86 |
| BEPAL w/o Auxiliary Learning | 8.35% | 39.37 |

TABLE III: Transferability study in PP Environment: Comparing native and transfer models for number of steps to complete the game

| | N=5 m=12 Max Steps = 40 obstacles = 0 | N=5, m=12 Max Steps = 40 obstacles = 10 | N=5, m=20 Max Steps = 80 obstacles = 20 |
|---|---|---|---|
| Native | 18.23 | 22.41 | 55.23 |
| Transfer | 19.76 | 22.41 | 54.62 |

73.80% reduction in standard error compared to the best baseline.

Notably, our method exhibits slower initial learning during the early epochs due to the additional loss introduced by auxiliary learning, which must be optimized concurrently. As a result, the learning process progresses more slowly at the beginning. However, once the auxiliary tasks are properly learned, our method surpasses all baselines in learning speed. This highlights the effectiveness of auxiliary learning in facilitating policy learning and improving overall performance.

It's important to point out that, although the standard deviation of the baselines decreases on the large map, this does mean their performance became more robust. Instead, this is due to the fact that agents using the baseline model often fail to complete the game within the maximum allowable steps. As a result, the number of steps taken remains to be the maximum value in most games.

*2) Ablation Study:* To quantify the impact of the auxiliary learning, we conducted an ablation study in an environment with a $12 \times 12$ map, 5 agents and 10 obstacles. Table II shows that after removing both reward and motion predictions, the number of steps taken by BEPAL-MAS increased 77.2%. When only the reward prediction or motion prediction was removed, the average number of steps increased by 10.4% and 66.0% respectively. The results support our hypothesis that performing auxiliary learning helps agents develop a robust understanding of the environment and achieves a better hidden state representations, and consequently, lead to a better game performance.

*3) Transferability:* To evaluate the transferability of BEPAL, we trained it in an environment with $12 \times 12$ map, 5 agents and 10 obstacles. The trained model is then transferred

to different game settings and its performance is tested. As a comparison, we also trained a native model using the configuration of the testing environment. The performance of the "Transferred" and "Native" models are compared in Table III. The results show that BEPAL has a robust performance and can be transferred to different testing environments.

It is interesting to note that, for large map (i.e., $m = 20$), the transfer model works better than the native model. This is probably because, given a more complex environment, the native trained model requires much longer training time. Hence, under the same number of training epoch, the transfer model performs better, which is an example of curriculum learning [31].

*4) Correlation between Auxiliary and Reinforcement Learning:* Experiments were conducted to demonstrate the correlation between the auxiliary learning and the MARL learning. Figure 3 illustrates the relationship between the accuracy of auxiliary prediction and the number of steps agents needed to complete the game. On average, the agent's performance in the game shows an average Pearson correlation coefficient of 0.41 with the accuracy of motion prediction and 0.55 with the accuracy of reward prediction on two games, both with a p-value close to 0. This positive correlation suggests that training the agent to perform auxiliary predictions enhances its performance in the game.

While the accuracy of reward prediction shows a higher correlation with agent performance than that of motion prediction, Table II shows that it has a comparatively smaller influence on improving agent performance. This can be attributed to the intrinsic relationship between reward prediction and the critic model. Since both the critic and the reward prediction model

(a) PP Motion Prediction      (b) PP Reward Prediction

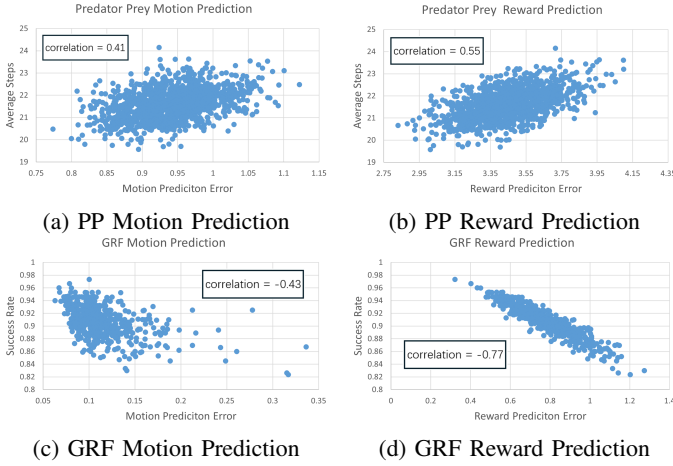(c) GRF Motion Prediction      (d) GRF Reward Prediction

Fig. 3: Game Performance vs. Auxiliary Prediction Accuracy.

are responsible for predicting rewards, albeit for different scopes (individual versus team), the performance of these two models are strongly correlated. Meanwhile, the quality of the critic model determines the success rate of the agent performance. Hence the accuracy of the reward predictor is strongly correlated to the game performance of the agent. Higher correlation with the agent performance does not mean the reward prediction provide more information to improve the RL training. On the contrary, because its high correlation with the the critic model, reward prediction does not have as significant impact to the RL training as motion prediction.

*5) Computation Overhead of Auxiliary Learning:* To evaluate whether the auxiliary tasks introduce significant computational overhead, we conducted experiments comparing the training times for models with and without auxiliary tasks. The experiments were performed on the same Intel(R) Core(TM) i7-8700K CPU @ 3.70GHz, with the average time per epoch calculated over 100 epochs. The BEPAL model with auxiliary tasks required an average of 73.98 seconds per epoch, while the model without auxiliary tasks took 68.67 seconds per epoch. This result demonstrates that incorporating auxiliary learning increases computation time by only 7.82% per epoch, while providing an average 16% improvement in performance. A learned auxiliary predictor also helps the RL training to converge faster, and eventually requires fewer epochs. These findings confirm that the addition of auxiliary tasks is computationally efficient and significantly beneficial to overall performance.

*6) Visualization of Auxiliary Predictions:* We visualize some motion predictions captured from a PP game. Fig. 4a-c show the motion prediction made by agent 3 in steps 11, 20, and 21; and Fig. 4d shows the prediction made by agent 1 in step 20. The prediction includes the coordinates and moving directions (represented as vectors) of all agents, including the predictor itself, in the next several time steps and the location of the prey. The prediction represents the agent's interpretation of the current environment state.

As shown in Fig. 4a, at time step 11 when no agent has

discovered the prey, agent 3 predicts that the prey location is in an unexplored area to motivate exploration. Its prediction of other agent's motion is a mixture of information from the received messages and its own prediction of the prey location. At time step 20, agent 1 observed the prey (Fig. 4d) and broadcast the information. Meanwhile, as shown in Fig. 4b, agent 3 and other agents have not yet received the message, they move from right side of the map, which has been explored, to the left side, which is unexplored. Finally, in step 21 (Fig. 4c), all agents received the message from agent 1, and they all move toward the true location of the prey. As we can see the predicted motion generally aligns with the actual motion of the agents. Without receiving new messages, the predicted movement distance (arrow length) for each agent (Fig. 4b) reduces as the agent is no long sure about the motion of its teammates. Once the agent acquires critical information, it starts predicting larger movement distances (Fig. 4c), reflecting increased certainty in its updated information.

Fig. 5 illustrates the reward predictions generated by agent 3 in the previous game. In the early steps of the game, the agent predicts future rewards with a negative average value, reflecting its learned expectation. As exploration progresses, the predicted rewards gradually increase, indicating steady progress in exploration. A sudden change in reward predictions signifies that the prey has been located. Subsequently, the agent updates its predictions of other agents' rewards, reflecting the belief that the game will soon end.

### D. Main Results for Google Research Football

The GRF environment features stochastic state transitions, sparse rewards, and a complex action space. The evaluation metric is the scoring success rate. We modified the GRF, originally a fully observable environment, into a partially observable one by masking objects outside the agent's visual range. It is important to note that after this modification, the performance of MAGIC and TARMAC declined significantly compared to previously reported results [7].

Figure 6 compares the change in the success rate over training epochs between BEPAL and three baseline models. Note that BAMS could not be applied to GRF due to its CNN-based observation encoder being incompatible with the non-grid environment of GRF. The results show that BEPAL converges to a higher success rate at a faster speed with more stable learning compared to all baselines.

Table II presents the results of the ablation study, demonstrating the impact of the two auxiliary learning tasks on the GRF game. As shown in the table, omitting both auxiliary tasks results in a 90% drop in the success rate. The success rate decreases by 12.7% and 77.1% when the reward prediction or motion prediction tasks are omitted, respectively. These results highlight the effectiveness of auxiliary learning in enhancing agent game performance.

The relationship between the prediction error and game performance is also shown in Figure 3. The effectiveness of both predictions follows the same trend as observed in the PP environment. However, we observe that GRF exhibits a
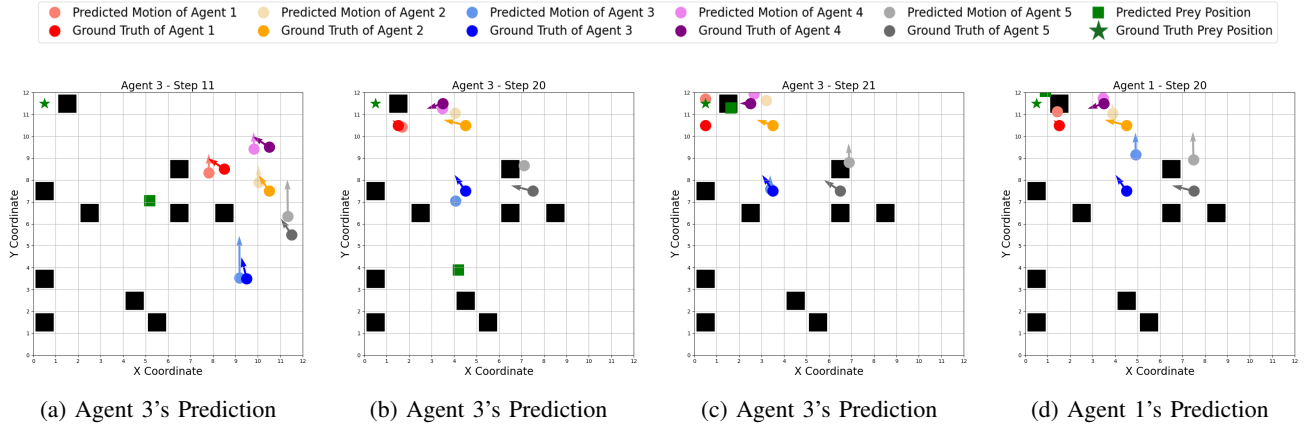
(a) Agent 3's Prediction    (b) Agent 3's Prediction    (c) Agent 3's Prediction    (d) Agent 1's Prediction

Fig. 4: Visualization of motion prediction generated by agents in different time steps in a $12 \times 12$ map with 10 obstacles. (a-c) Predictions generated by agent 3 at time steps 11, 20, and 21. (d) Prediction generated by agent 1 in time step 20.
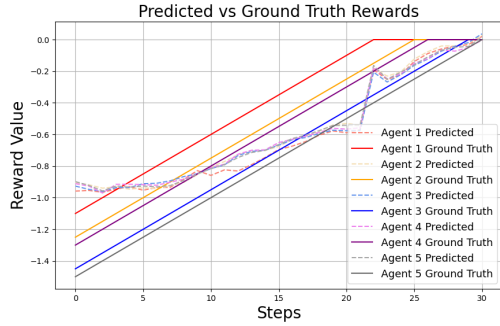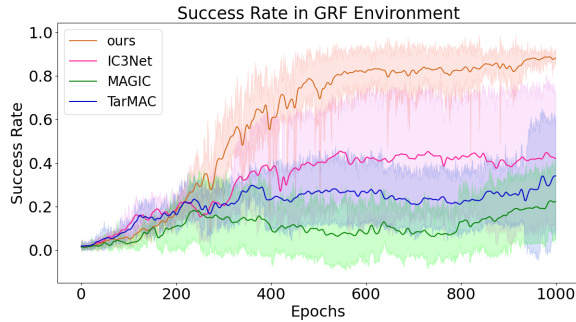


Fig. 5: Visualized Reward Prediction



Fig. 6: Comparison of baselines on GRF environment

particularly strong correlation between reward prediction and success rate. This is because GRF provides a global reward shared by all agents. The more accurate an agent can predict team member's future reward, the better it can predict its own expected reward, which consequently lead to better action. However, as discussed in Section V-C4, a stronger correlation with reward prediction does not necessarily indicate a greater influence on final performance.

## VI. CONCLUSIONS

This paper proposes a novel training approach that uses auxiliary learning to enhance multi-agent reinforcement learning in a partially observable environment. The auxiliary objective functions provide additional feedback during training, refining the agent's hidden state representation and enriching it with information that, while not directly contributing to the agent's reward, benefits in-agent collaboration. The auxiliary learning leverages a belief decoder to predict its teammates' motion and reward information from the agent's hidden state. By training on these auxiliary prediction tasks, the agent gains a better understanding of its environment and becomes more aware of its teammates' status.

We compared our method with IC3Net, MAGIC, and Tar-MAC in both predator-prey environments (with and without obstacles) and the Google Research Football (GRF) environment. Our approach demonstrates that auxiliary learning significantly improves the agent's game performance and makes the training process more robust.

## REFERENCES

[1] L. Kraemer and B. Banerjee, "Multi-agent reinforcement learning as a rehearsal for decentralized planning," *Neurocomputing*, vol. 190, pp. 82–94, 2016.

[2] F. A. Oliehoek, M. T. Spaan, and N. Vlassis, "Optimal and approximate q-value functions for decentralized pomdps," *Journal of Artificial Intelligence Research*, vol. 32, pp. 289–353, 2008.

[3] A. Oroojlooy and D. Hajinezhad, "A review of cooperative multi-agent deep reinforcement learning," *Applied Intelligence*, vol. 53, no. 11, pp. 13677–13722, 2023.

[4] D. Simões, N. Lau, and L. P. Reis, "Multi-agent actor centralized-critic with communication," *Neurocomputing*, vol. 390, pp. 40–56, 2020.

[5] S. Iqbal and F. Sha, "Actor-attention-critic for multi-agent reinforcement learning," in *International conference on machine learning*, pp. 2961–2970, PMLR, 2019.

[6] A. Singh, T. Jain, and S. Sukhbaatar, "Learning when to communicate at scale in multiagent cooperative and competitive tasks," *arXiv preprint arXiv:1812.09755*, 2018.

[7] Y. Niu, R. R. Paleja, and M. C. Gombolay, "Multi-agent graph-attention communication and teaming.," in *AAMAS*, vol. 21, p. 20th, 2021.

[8] A. Das, T. Gervet, J. Romoff, D. Batra, D. Parikh, M. Rabbat, and J. Pineau, "Tarmac: Targeted multi-agent communication," in *International Conference on machine learning*, pp. 1538–1546, PMLR, 2019.

[9] J. A. Ayala-Romero, P. Mernyei, B. Shi, and D. Mazón, "Kraf: A flexible advertising framework using knowledge graph-enriched multi-agent reinforcement learning," in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pp. 47–56, 2022.

[10] Z. Pu, H. Wang, Z. Liu, J. Yi, and S. Wu, "Attention enhanced reinforcement learning for multi agent cooperation," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[11] H. Lin, C. Lyu, Y. He, Y. Liu, K. Gao, and X. Qu, "Enhancing state representation in multi-agent reinforcement learning for platoon-following models," *IEEE Transactions on Vehicular Technology*, 2024.

[12] S. Sukhbaatar, R. Fergus, *et al.*, "Learning multiagent communication with backpropagation," *Advances in neural information processing systems*, vol. 29, 2016.

[13] C. Luo, Q. Huang, A. B. Wu, S. Khan, H. Li, and Q. Qiu, "Multi-agent cooperative games using belief map assisted training," in *ECAI 2023*, pp. 1617–1624, IOS Press, 2023.

[14] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.

[15] D. Ye, M. Zhang, and Y. Yang, "A multi-agent framework for packet routing in wireless sensor networks," *sensors*, vol. 15, no. 5, pp. 10026–10047, 2015.

[16] Z. Xu, Y. Lyu, Q. Pan, J. Hu, C. Zhao, and S. Liu, "Multi-vehicle flocking control with deep deterministic policy gradient method," in *2018 IEEE 14th International Conference on Control and Automation (ICCA)*, pp. 306–311, IEEE, 2018.

[17] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson, "Monotonic value function factorisation for deep multi-agent reinforcement learning," *Journal of Machine Learning Research*, vol. 21, no. 178, pp. 1–51, 2020.

[18] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *Advances in neural information processing systems*, vol. 30, 2017.

[19] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, 2018.

[20] X. Yao, C. Wen, Y. Wang, and X. Tan, "Smix ($\lambda$): Enhancing centralized value functions for cooperative multiagent reinforcement learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, no. 1, pp. 52–63, 2021.

[21] Y. Yang, J. Hao, G. Chen, H. Tang, Y. Chen, Y. Hu, C. Fan, and Z. Wei, "Q-value path decomposition for deep multiagent reinforcement learning," in *International Conference on Machine Learning*, pp. 10706–10715, PMLR, 2020.

[22] M. Zhou, Z. Liu, P. Sui, Y. Li, and Y. Y. Chung, "Learning implicit credit assignment for cooperative multi-agent reinforcement learning," *Advances in neural information processing systems*, vol. 33, pp. 11853–11864, 2020.

[23] J. Su, S. Adams, and P. Beling, "Value-decomposition multi-agent actor-critics," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, pp. 11352–11360, 2021.

[24] J. Foerster, I. A. Assael, N. De Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," *Advances in neural information processing systems*, vol. 29, 2016.

[25] Y. Wen, Y. Yang, R. Luo, J. Wang, and W. Pan, "Probabilistic recursive reasoning for multi-agent reinforcement learning," *arXiv preprint arXiv:1901.09207*, 2019.

[26] M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu, "Reinforcement learning with unsupervised auxiliary tasks," *arXiv preprint arXiv:1611.05397*, 2016.

[27] E. Shelhamer, P. Mahmoudieh, M. Argus, and T. Darrell, "Loss is its own reward: Self-supervision for reinforcement learning," *arXiv preprint arXiv:1612.07307*, 2016.

[28] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[29] G. Wayne, C.-C. Hung, D. Amos, M. Mirza, A. Ahuja, A. Grabska-Barwinska, J. Rae, P. Mirowski, J. Z. Leibo, A. Santoro, *et al.*, "Unsupervised predictive memory in a goal-directed agent," *arXiv preprint arXiv:1803.10760*, 2018.

[30] S. Ruder, "An overview of gradient descent optimization algorithms," *arXiv preprint arXiv:1609.04747*, 2016.

[31] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proceedings of the 26th annual international conference on machine learning*, pp. 41–48, 2009.

[32] Y. Liu, W. Wang, Y. Hu, J. Hao, X. Chen, and Y. Gao, "Multi-agent game abstraction via graph attention neural network," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, pp. 7211–7218, 2020.

[33] K. Kurach, A. Raichuk, P. Stańczyk, M. Zając, O. Bachem, L. Espeholt, C. Riquelme, D. Vincent, M. Michalski, O. Bousquet, *et al.*, "Google research football: A novel reinforcement learning environment," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, pp. 4501–4510, 2020.

[34] Z. Zhao, F. Zhao, Y. Zhao, Y. Zeng, and Y. Sun, "A brain-inspired theory of mind spiking neural network improves multi-agent cooperation and competition," *Patterns*, vol. 4, no. 8, 2023.

[35] C. L. Sebastian, N. M. Fontaine, G. Bird, S.-J. Blakemore, S. A. De Brito, E. J. McCrory, and E. Viding, "Neural processing associated with cognitive and affective theory of mind in adolescents and adults," *Social cognitive and affective neuroscience*, vol. 7, no. 1, pp. 53–63, 2012.

[36] C. Fang and K. L. Stachenfeld, "Predictive auxiliary objectives in deep rl mimic learning in the brain," *arXiv preprint arXiv:2310.06089*, 2023.

[37] R. Huang, X. Chu, J. Zhang, and Y. H. Hu, "Energy-efficient monitoring in software defined wireless sensor networks using reinforcement learning: A prototype," *International Journal of Distributed Sensor Networks*, vol. 11, no. 10, p. 360428, 2015.

[38] K. O'shea and R. Nash, "An introduction to convolutional neural networks," *arXiv preprint arXiv:1511.08458*, 2015.