

# Residual Q-Networks for Value Function Factorizing in Multi-Agent Reinforcement Learning

Rafael Pina, Varuna De Silva, Joosep Hook and Ahmet Kondo, *Senior Member, IEEE*

**Abstract**—Multi-Agent Reinforcement Learning (MARL) is useful in many problems that require the cooperation and coordination of multiple agents. Learning optimal policies using reinforcement learning in a multi-agent setting can be very difficult as the number of agents increases. Recent solutions such as Value Decomposition Networks (VDN), QMIX, QTRAN and QPLEX adhere to the centralized training and decentralized execution scheme and perform factorization of the joint action-value functions. However, these methods still suffer from increased environmental complexity, and at times fail to converge in a stable manner. We propose a novel concept of Residual Q-Networks (RQNs) for MARL, which learns to transform the individual Q-value trajectories in a way that preserves the Individual-Global-Max criteria (IGM), but is more robust in factorizing action-value functions. The RQN acts as an auxiliary network that accelerates convergence and will become obsolete as the agents reach the training objectives. The performance of the proposed method is compared against several state-of-the-art techniques such as QPLEX, QMIX, QTRAN and VDN, in a range of multi-agent cooperative tasks. The results illustrate that the proposed method, in general, converges faster, with increased stability and shows robust performance in a wider family of environments. The improvements in results are more prominent in environments with severe punishments for non-cooperative behaviours and especially in the absence of complete state information during training time.

**Index Terms**—Multi-agent reinforcement learning (MARL), value function factorization, deep learning, task cooperation.

## I. INTRODUCTION

REINFORCEMENT learning has been growing as one of the most prominent areas of artificial intelligence in the recent past [1]. It has proved to be successful in single-agent tasks in diverse areas of application such as finance, games, and manufacturing [2], [3]. There are many applications that require reinforcement learning to be done in Multi-Agent settings, which require some level of cooperation and coordination among the agents [4].

In centralized Multi-Agent Reinforcement Learning (MARL), observations and actions of all agents are considered as a whole by a centralized oracle. However, centralised learning suffers from Curse of Dimensionality as the number of agents increases and the action space becomes extremely large [5], [6]. On the other end of the spectrum, fully decentralized MARL treats each agent as an independent learner. However, learning in such a setting becomes non-stationary from the perspective of a single agent, and renders it difficult for an agent to understand the reasons for the joint reward. As a result, the agent might get confounded by the rewards of other agents. Thus, the optimal policy may not be learnable in such settings. This problem is

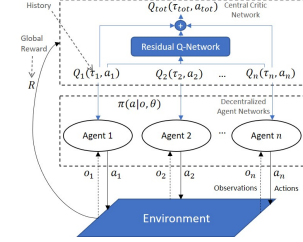


Fig. 1: Illustration of the centralized training decentralized execution paradigm within the proposed RQN method.

referred to as relative overgeneralization [7]. To minimize the above problems, a new paradigm has been investigated for MARL, known as the centralized training and decentralized execution (CTDE). It allows the centralized critic to use the full information available of the environment during training, but restricts agents (actors) to their local observations when it comes to execution. Most of the MARL methods proposed in the recent past focus on the centralized training decentralized execution scheme [8]–[11]. This paradigm helps to minimize the relative overgeneralization problem by improving the credit assignment to the agents, which helps agents to appropriately deduce their individual contribution to the overall team reward [9]. Yet, a fully centralized critic still suffers from Curse of Dimensionality [12].

One possible approach to overcome the Curse of Dimensionality is to learn a centralized, but a factorized critic, which learns to decompose the joint action-value function to individual action-value functions of the agents. The Value Decomposition Networks (VDN) [13] and QMIX [12], are two examples of such techniques. VDN aims to learn an additive decomposition of the joint action-value function. However, authors in [12] argue that the simple linear operations performed by VDN limit the complexity of the value functions that it can represent. QMIX tries to extend the range of value functions that can be represented by mixing the individual Q-values in a monotonic way. However, VDN and QMIX are limited by additivity and monotonicity constraints, respectively, and are unable to factorize action-value functions that do not follow these constraints. [14] proposed QTRAN, a solution capable of overcoming these problems by using a better factorization process. The main idea of QTRAN is to transform the original action-value function into a new one, and factorize the transformed function. Despite the many solutions brought by these methods, there are still problems to tackle that have been mitigated but not totally overcome yet such as the credit assignment problem [15].

In this paper we show that there exists a more relaxed condition than in QTRAN [14] under which value function factorization is possible. Accordingly, we propose a new deep learning based method that aims to learn factorizations of the action-value functions in different tasks by performing complex non-linear operations over the individual Q-values. We propose the use of a new Residual Q-Network (RQN) (Figure 1) that will allow to reduce the impact of the credit assignment problem [15], [16]. In order to achieve that, the RQN will learn an individual factor for each agent that will indicate the relative importance of a certain Q-value trajectory (the sequence of action-values of a given episode). As a result, the proposed method is able to also improve the credit assignment of the agents and achieve a quick and stable level of optimal performance over time by learning estimations for the individual Q-values. The main role of the RQN is to learn estimations for the individual Q-values in order to coax the agents towards optimal performances. The intuition is that if individual agents are performing well by visiting certain states with relatively high Q-values, they should be encouraged to visit those states in subsequent episodes. To do that, the proposed model should take advantage of the individual Q-value trajectories. The empirical results suggest that the proposed approach is affected less by increasing the number of agents, when compared to previous value decomposition-based methods. The performance of the proposed method is benchmarked against previous related approaches (QMIX [12], QTRAN [14], VDN [13], QPLEX [17] and WQMIX [18]) in a set of cooperative multi-agent environments of varying complexity. Contributions of this paper are as follows:

- 1) We present novel theoretical evidence for factorization conditions in a wider family of different environments.
- 2) We propose a new method for value function factorization in Multi-Agent Reinforcement Learning cooperative tasks that uses a novel neural network architecture based on a novel concept of Residual Q-Networks (RQNs).
- 3) We demonstrate the viability of the proposed method in a diverse set of environments with varying levels of complexity, with distinct tasks and number of agents.

The rest of the paper is organized as follows: In section II we present recent literature that is associated with the contributions of this paper, and section III presents the theoretical details fundamental to the developments of this paper. Section IV presents the proposed novel value factorization methodology for MARL tasks. In section V, experimental details are presented, followed by results and discussions in section VI. We conclude the paper in section VII with indications for future work.

## II. RELATED WORK

[19] introduced one of the first Multi-Agent Reinforcement Learning approaches, based on independent Q-learning [20]. In a more complex but related approach, [21] extended the concept with Deep Q-Networks to multi-agent environments, where individual agents controlled by Independent Deep Q-Networks learn both cooperative and competitive behaviours.

Deep reinforcement learning for multi-agent settings quickly became popular, but the increase of the number of

agents proved to be a difficult problem to solve. To mitigate this problem, a recurring approach in recent works is the use of a centralized training and decentralized execution model [5]. COMA [9] is an example of an actor-critic centralized training decentralized execution method that uses a centralized critic to estimate values for state-action pairs and a decentralized actor that maps states to actions. Other recent works also follow this paradigm such as [8], [16] or [18]. Methods such as QMIX [12], VDN [13], QTRAN [14] and QPLEX [17] also adhere to the centralized training and decentralized execution approach. Moreover, they aim to learn a factorization for the centralized critic, which proved to be a successful strategy. Further, these methods employ a parameter sharing approach to improve learning efficiency [10]. However, VDN and QMIX are constrained by additivity and monotonicity, respectively [14]. In addition, methods like QMIX, QPLEX or QTRAN take advantage of the global state of the environment, which can be unfeasible under certain conditions. On the side, works such as [22] use similar approaches together with an idea of task decomposition to improve the learning performance.

Communication-based approaches are also popular in multi-agent scenarios [2], [8], [23], [24]. A key difference of this family of algorithms from fully decentralized execution, is that the agents can communicate among themselves during the execution process. Multiple different works proved that this model can be effective. [25] proposes an algorithm where the agents learn to coordinate by performing multiple rounds of communication among themselves before taking an action in the environment. Also, [23] demonstrate interesting results by showing that a group of agents is capable of learning complex communication protocols during centralized training phase, to solve cooperative tasks with communication during execution time.

The proposed work in this paper is a new method that follows a centralized training and decentralized execution model for cooperative scenarios. Accordingly, it falls in the category of VDN, QMIX, QTRAN and QPLEX where an effective decomposition of the joint action-value function is learnt. In section III we elaborate on the theoretical details of these methods.

## III. BACKGROUND

We consider a set of fully cooperative tasks that can be represented as Decentralized Partially Observable Markov Decision Processes (Dec-POMDP) [26]. A Dec-POMDP is an extension of a Markov Decision Process to a multi-agent setting where the environment is not fully perceptible by each agent. Let the tuple  $G = \langle S, A, R, O, Z, P, N, \gamma \rangle$  be the representation of a Dec-POMDP, where  $S$  is a set of states  $s$  that represent the current state of the environment. For each agent  $i \in \mathcal{N} \equiv \{1, \dots, N\}$  consider  $A$  to be the set of actions  $a_i$  that each agent can choose at a given time step  $t$ ,  $a_i \in A : A = \{a_1, \dots, a_n\}$  and  $R(s, a) : S \times A \rightarrow \mathbb{R}$  and  $O(s, i) : S \times \mathcal{N} \rightarrow \mathcal{Z}$  are the reward and observation functions, respectively. The reward function is shared by all the agents and  $Z$  represents the model of observations. The discount factor of the reward function is represented by  $\gamma$  :

$\gamma \in [0, 1]$ . Since we consider a partially observable system, each agent receives a local observation  $o_i$  from the model  $O$ . Furthermore, each agent  $i$  holds an action-observation history represented by  $\tau_i$  so that  $\tau_i \in T \equiv (Z \times A)^* \rightarrow \{\tau_1, \dots, \tau_N\}$ . A transition of state in the environment is performed according to the probability function  $P(s'|s, a) : S \times A \times S \rightarrow [0, 1]$ , where  $s'$  represents the state that follows  $s$  after taking an action  $a$ . This constructs a certain policy  $\pi_i(a_i|\tau_i)$ . The objective is to find an optimal joint policy  $\pi$  that maximizes the joint action-value function  $Q_\pi(s_t, a_t) = \mathbb{E}_\pi[R_t|s_t, a_t]$ , where  $R_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k}$  is the discounted return.

In this section, we introduce the key concepts to understand the method proposed in this paper.

#### A. Reinforcement Learning with Deep Q-Networks

In the work of [27], Deep Q-Networks (DQN) is an approach that uses deep neural networks to approximate value-functions  $Q(s, a)$ . This is done using a target network, whose parameters are copied every  $n$  steps from an evaluation network, where  $n$  is a pre-defined number of steps to update the target network. Furthermore, the networks have access to a buffer with an observation-action history used to train them. In Deep Q-Learning the loss of the networks is calculated according to

$$\mathcal{L}(\theta) = [y^{DQN} - Q(s, a; \theta)]^2 \quad (1)$$

where  $y^{DQN}$  is the DQN target network,  $y^{DQN} = R(s, a) + \gamma \max_{a'} Q(s', a'; \theta^-)$ , and  $\theta$  and  $\theta^-$  represent the parameters of the evaluation network and target, respectively.

#### B. Value Function Factorization in Centralized Training-Decentralized Execution Based MARL

Under the centralized training decentralized execution paradigm, recent studies defined a key condition that should be satisfied to enable efficient learning during the centralized training phase. This condition is called Individual-Global-Max (IGM) [14] and is defined as

$$\arg\max_a Q_{tot}(\tau, a) = \begin{pmatrix} \arg\max_{a_1} Q_1(\tau_1, a_1) \\ \vdots \\ \arg\max_{a_N} Q_N(\tau_N, a_N) \end{pmatrix} \quad (2)$$

In simple terms, we can describe this property as requiring the optimal joint action to be the same as the combination of locally optimal actions of individual agents.

Value Decomposition Networks [13] is a method that aims to learn an additive decomposition over individual agents based on the team reward. The key idea is the assumption that a joint action-value function of a group of agents can be represented as the addition of individual action-value functions of each agent. In simple terms, this method decomposes the action-value function according to

$$Q_{tot}(\tau, a) = \sum_{i=1}^N Q_i(\tau_i, a_i; \theta_i) \quad (3)$$

where  $Q_{tot}$  is the joint action-value function and  $Q_i$  represents the individual value functions for each agent  $i$  with parameters

$\theta_i$ . However, it is argued that this type of decomposition is too simple to represent a wide range of environments due to the additivity constraint [12], [14].

QMIX is a method that learns a more complex decomposition of the joint action-value function using hyper networks that take advantage of extra state information during centralized training [12]. The decomposition is done using a mixing network that uses complex non-linear operations that follow an imposed monotonicity constraint, represented as

$$\frac{\partial Q_{tot}}{\partial Q_i} \geq 0, \forall i \quad (4)$$

The motivation for this monotonicity constraint is to generalize the family of functions that VDN can represent to a larger family of monotonic functions. This condition is assured by keeping the weights of the mixing network positive. This constraint has proved to be a problem when it comes to solving tasks that do not meet the monotonicity conditions [14], [18]. For example, a task where the best action of one agent depends on the actions of a second agent, at the same time step. We demonstrate this effect using a matrix game in subsection IV-C. To update their networks, both VDN and QMIX use the DQN loss as stated in (1), but their losses operate on  $Q_{tot}$  instead of an individual  $Q$  of DQN.

QTRAN [14] is a different value-decomposition approach that proposes a new method of factorizing joint action-value functions without the constraints imposed by VDN and QMIX: additivity and monotonicity. The key idea of QTRAN is to transform the original action-value function using an additive decomposition, and then to correct it by using a joint state-dependent value function. To address the additivity and monotonicity constraints present in VDN and QMIX, QTRAN proposes a new approach that uses multiple networks that are trained in a centralized way. QTRAN uses an individual action-value network for each agent  $i$ ,  $f_q: (\tau_i, a_i) \mapsto Q_i$ , a joint action-value network,  $f_r: (\tau, a) \mapsto Q_{tot}$ , and a state-value network,  $f_v: \tau \mapsto V_{tot}$ . In this sense, QTRAN uses a weighted loss that aims to both learn the optimal  $Q_{tot}$  and factorize the transformed action-value function. The use of an unconstrained joint action-value function is key to keep a wide representational complexity. With the method proposed in this paper, we also intend to eliminate both the additivity and the monotonicity constraints by adding more layers of complexity to the decomposition procedure. This process will be explained in the sections ahead.

#### C. Value Mapping in Residual Networks

ResNets was proposed by [28] as a novel approach for image recognition. The key to these networks is the existence of skip connections that allow values to skip some layers in the network. Similarly the work in this paper introduces a novel architecture for MARL with the existence of skip connections. The works of [28] show the strength of such architecture that proves to be easier to optimize than previous networks.

In a Residual Network, the computations performed to map subsequent values are obtained through the equation [29]

$$x_U = x_u + \sum_{i=u}^{U-1} F(x_i, w_i) \quad (5)$$

where  $F$  is a non-linear function with inputs  $x$ , parameterized by  $w$ , and  $x_U$  are the features of a deeper unit  $U$  and  $x_u$  are the features of a shallower unit  $u$ .

#### IV. RESIDUAL Q-NETWORK (RQN) FOR VALUE FUNCTION FACTORIZATION

The proposed method involves a centralized, but a factorized critic network supplemented by a novel Residual Q-Network (RQN). In this section we show a more relaxed factorization condition that suits a wider family of joint action-value functions. With the proposed new network architecture, we aim to exploit this condition.

##### A. Requirements for Factorization with Individual Factor Functions

We start by defining a Theorem over the Theorem 1 of value function factorization in [14], which defines sufficient conditions to satisfy IGM in (2).

**Theorem 1.** *Given a joint action-value function  $Q_{tot}(\tau, a)$ , we say that it is factorized by  $[Q_i(\tau_i, a_i)]$  if*

$$\sum_{i=1}^N (Q_i(\tau_i, a_i) + \phi_i(\tau)) - Q_{tot}(\tau, a) = \begin{cases} 0 & a = \bar{a}, \\ \geq 0 & a \neq \bar{a}, \end{cases} \quad (6a) \quad (6b)$$

where  $\phi_i$  are the individual correction factors and  $\sum_{i=1}^N \phi_i(\tau) = \max_a Q_{tot}(\tau, a) - \sum_{i=1}^N Q_i(\tau_i, \bar{a}_i)$ .

Accordingly, we utilize a transformation of individual action-value functions that preserve IGM,  $\sum_{i=1}^N (Q_i(\tau_i, a_i) + \phi_i(\tau))$ . This is an affine transformation of the individual factor functions, under which the IGM condition is not violated. This is a much stricter condition than additivity constraint of VDN and monotonicity constraint of QMIX [14]. Note that this transformation also shares some theoretical similarities with the individual dueling transformation introduced in QPLEX [17]. However, RQN learns a trajectory-based factor and uses this factor to transform the individual  $Q_i$  while in the transformation of QPLEX, the authors use advantage functions together with the state-value functions.

*Proof.* Theorem 1 shows that if the above condition holds, then  $Q_i$  satisfies IGM. In a similar manner to [14], we show that if there is a  $Q_i$  that satisfies the conditions above, then the joint action that maximizes the  $Q_{tot}$  is the same as the set of optimal local actions  $\bar{a} = [\bar{a}_i]_{i=1}^N$ , where  $a_i$  denotes an optimal local action of the agent  $i$ ,  $\bar{a}_i = \arg\max_{a_i} Q_i(\tau_i, a_i)$ . We have that

$$\begin{aligned} Q_{tot}(\tau, \bar{a}) &= \sum_{i=1}^N (Q_i(\tau_i, \bar{a}_i) + \phi_i(\tau)) && \text{from (6a)} \\ &= \sum_{i=1}^N (Q_i(\tau_i, \bar{a}_i)) + \sum_{i=1}^N \phi_i(\tau) \\ &\geq \sum_{i=1}^N (Q_i(\tau_i, a_i)) + \sum_{i=1}^N \phi_i(\tau) \\ &\geq Q_{tot}(\tau, a) && \text{from (6b)} \end{aligned}$$

This means that the set of optimal local actions  $\bar{a}$  maximizes  $Q_{tot}$ , i.e.,  $Q_i$  satisfies IGM. End of Proof.  $\square$

##### B. Factorizing with Residual Q-Networks

In this section, we propose a deep reinforcement learning method that exploits Theorem 1. Our motivation is to achieve an effective decomposition by learning a new estimation factor that will affine transform the individual Q-value trajectories  $Q_i(\tau_i, a_i)$ . A different factor  $\phi_i(\tau)$  is learnt for each agent, so that we can prioritize different state-action trajectories taken by agents within an episode.

The proposed architecture for the centralized critic is illustrated in Figure 2. The individual agent networks are recurrent neural networks using a GRU with width 64. These networks receive the observations  $o_i$  from the agents at a given time step and will output individual action-values based on the observation-action histories. These values are passed on to the RQN, where a set of non-linear operations are performed. Over a batch of episodes, an estimation network (within the RQN) will take as input features the mean of the Q-values of each agent over the steps of each episode, and their maximum Q-value. This estimation network has an input shape of  $2 \times N$ , where  $N$  is the number of agents. The inputs are given to a Linear layer with 64 units, followed by a ReLU non-linearity and finally another Linear layer, where the output has size equal to  $N$ . In other words, the output is a set with  $N$  estimation factors, one for each agent. The estimation factors are added to the initial action-values from the individual agent networks, to compute the adjusted individual Q-values. The adjusted Q-values are given as inputs to the summing network, to calculate the joint action-value function.

We denote estimation factor, by which the individual Q-values are adjusted, as  $\phi_i$  for each agent  $i : i \in \{1, \dots, N\}$ . This factor is the output of the RQN. Learning this factor successfully is the first aspect of the proposed method. The network is trained by receiving two features for each agent, one being the mean of the Q-values across all the steps  $t : t \in \{1, \dots, T\}$  of each episode  $e$ ,

$$\bar{x}_i^e = \frac{\sum_{t=0}^T Q_i^t(s_i^t, a_i^t)}{T} \quad (7)$$

and the second is the maximum Q-value of each agent across all the steps of each episode  $e$ ,

$$m_i^e = \max_t Q_i^t(s_i^t, a_i^t) \quad (8)$$

The intuition to include the maximum as a feature to the estimation network (see Figure 2) is that if the trajectory of an agent contains a relatively high value state-action pair, it should be encouraged to explore it in subsequent episodes. Further, mean value of a trajectory signals the relative goodness of an entire trajectory and maximum value will signify the goodness of some points within the trajectory. The estimation factor for each agent and an episode  $e$  can be calculated according to

$$\phi_i = f([\bar{x}_i^e], [m_i^e]; \theta)_i, \quad i \in \{1, \dots, N\} \quad (9)$$

where  $f$  represents the RQN with parameters  $\theta$ . This estimation factor, which is a non-linear function of the action-value trajectories, allows to decompose action-value functions of more complex environments than VDN due to the elimination of the simple sum constraint [12]. Furthermore, the



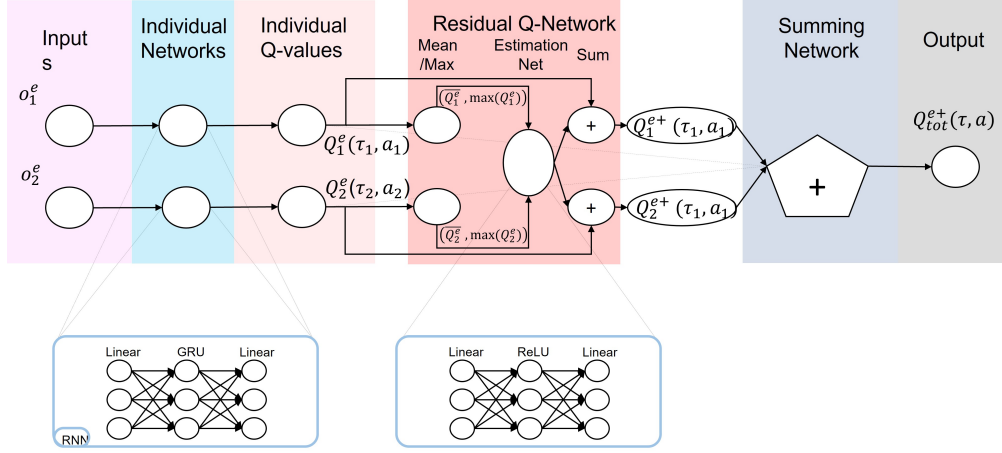


Fig. 2: Architecture of the centralized critic where the RQN is described in detail considering an example with two agents, 1 and 2, and for a certain episode  $e$ .

factorization procedure of the proposed method can benefit from the trajectories of the individual Q-values. In contrast to QMIX, the RQN based method removes the monotonicity constraint by not enforcing the weights of the network to remain positive. In comparison to QTRAN, where the state-value network computes a single value for  $Q_{tot}$ , the RQN calculates an estimation factor for each individual Q-values separately.

After the estimation factor for each agent is calculated by the network, the factor is added to the respective Q-value trajectories resulting in a new estimated Q-value for each agent  $i$ , according to  $Q_i^+ = Q_i + \phi_i$ , where  $\phi_i$  is obtained from (9).

At this point, we have learned an estimated version of the Q-values for each agent, the estimated Q-values, denoted by  $Q_i^+$ . Now, the critic network can be used to compute a new joint action-value function using the estimated Q-values,

$$Q_{tot}^+ = \sum_{i=1}^N Q_i^+ (\tau_i, a_i; \theta_i) \quad (10)$$

The proposed method is capable of learning a decomposition more capable of factorizing different factorizable tasks than previous methods due to the use of a RQN. As described in (5) consider now that  $H(Q_i) = F(x) + Q_i$  gives the output of the RQN. Recalling (7) and (8), consider that  $x = (\bar{x}_i^e, m_i^e)$  is a set containing the inputs for an agent  $i$  and an episode  $e$  to the estimation network of our model (Figure 2). Consider the example where the optimal output of the RQN given this input  $Q_i$  would be  $H(Q_i) = Q_i$ . This would be difficult to approximate using the normal mapping of a common neural network [28], since by inputting  $Q_i$  then it would need to be optimally mapped to  $Q_i$  after going through the network. On the other hand, if we can compute a value  $F(x)$  given the input  $x$  as a combination of features regarding the original input to the RQN, and  $F$  being the estimation network, then we end up with a factor  $\phi_i$  (as represented in (9)) given by  $F(x) = \phi_i$  that will converge to a stable value when the training objective is reached, since at that point, the role of the RQN as an auxiliary network is no longer needed. We investigate this effect further in subsection VI-D. Thus,  $H(Q_i) = F(x) + Q_i$

will be easier to learn with the RQN than with a common neural network since  $F(x)$  stabilizes and  $Q_i$  skips connections as described in III-C.

### C. Comparison of Representational Capacity of RQN with QTRAN, VDN and QMIX

For this demonstration, let  $Q_{jt}$  represent a joint action-value function. The key idea of QTRAN is to transform the original joint action-value function  $Q_{jt}$  into a new one  $Q'_{jt}$  that shares the optimal joint action with  $Q_{jt}$ . Theorem 1 in QTRAN [14] states a sufficient condition for  $[Q_i]$  that satisfies IGM,  $i \in \{1, \dots, N\}$ . Accordingly, a factorizable joint action value function  $Q_{jt}(\tau, a)$  is factorized by  $[Q_i(\tau_i, a_i)]$ , if

$$\sum_{i=1}^N Q_i(\tau_i, a_i) - Q_{jt}(\tau, a) + V_{jt}(\tau) = \begin{cases} 0 & a = \bar{a}, \\ \geq 0 & a \neq \bar{a}, \end{cases} \quad (11a)$$

where  $V_{jt} = \max_a Q_{jt}(\tau, a) - \sum_{i=1}^N Q_i(\tau_i, \bar{a}_i)$  and  $\bar{a}$  is a set of optimal local actions  $\bar{a}_i$ .

The above condition becomes necessary under an affine transformation of action-value functions,  $\psi = A \cdot Q + B$ , where  $A = [a_{ii}] \in R_+^{N \times N}$  is a symmetric diagonal matrix and  $B = [b_i] \in R^N$ . An important observation is that values of  $B$  does not alter the IGM condition. Thus, QTRAN relies on a linear transformation of individual factor functions,  $Q'_{jt}(\tau, a) = \sum_{i=1}^N Q_i(\tau_i, a_i)$ , a transformation that contains IGM. Thus, the transformation used in QTRAN is an identity matrix.

The representational power of QTRAN over VDN relies on the  $V_{jt}(\tau)$ , which tries to learn the difference (or the residual) between the joint action-value function  $Q_{jt}(\tau, a)$  and the transformation at individually optimal actions  $Q'_{jt}(\tau, \bar{a})$ . Thus  $V_{jt}(\tau)$ , is a residual function, that acts on top of the transformation.

In comparison, the proposed RQN, relies on an affine transformation, on individual factor functions. Thus, giving more flexibility for learning the individual maximums, while

retaining the IGM property beyond additivity and monotonicity constraints.

In comparison to QTRAN, the transformation adopted by RQN is as follows, where  $\phi_i$  are individual correction factors:

$$Q'_{jt} = \sum_{i=1}^N (Q_i(\tau_i, a_i) + \phi_i(\tau)) \quad (12)$$

In comparison to QTRAN, the representational difference with VDN, in RQN comes from  $\sum_{i=1}^N \phi_i$ .

The success of the proposed RQN structure is due to the difference in calculating the residual factors. In QTRAN, the authors rely on the joint action value function  $Q_{tot}(\tau, a)$ , to learn the residual between transformation, whereas in RQN we rely on a transformation that adheres to the sufficient conditions of IGM as much as QTRAN, but the learning problem is formulated to finding the individual maximums.

To further illustrate the differences in the representational complexities of the different approaches, we show the reconstructions of the  $Q_{tot}$  in a one-step matrix game, as used in [14] and [18]. In this simple game, two agents can choose one of three different actions in order to optimize a joint action, given by  $Q_{tot}$ . The matrix of the game is a non-monotonic matrix represented by Table Ia. The agents are trained for 20000 episodes (1-step episodes) in full exploration ( $\epsilon = 1$ ) and with a replay buffer of size 500. Tables Ib-Ie show the reconstructed  $Q_{tot}$  for VDN, QMIX, QTRAN and RQN, respectively. As the tables show, only QTRAN and RQN are capable of successfully reconstructing this non-monotonic matrix, whereas VDN and QMIX fail. This supports the fact that RQN can represent a wider family of action-value functions, eliminating constraints such as monotonicity or additivity.

TABLE I: Payoff matrix of the one-step game and reconstructed  $Q_{tot}$  for VDN, QMIX, QTRAN, and RQN.

$a_0 \backslash a_1$	A	B	C
A	8	-12	-12
B	-12	0	0
C	-12	0	0

(a) Payoff matrix

$a_0 \backslash a_1$	A	B	C
A	-7.50	-5.72	-5.18
B	-5.70	-3.91	-3.37
C	-4.97	-3.19	-2.65

(b) VDN

$a_0 \backslash a_1$	A	B	C
A	-7.55	-7.55	-7.55
B	-7.55	-0.01	-0.01
C	-7.55	-0.01	-0.01

(c) QMIX

$a_0 \backslash a_1$	A	B	C
A	7.97	-12.00	-11.99
B	-12.00	0.00	0.00
C	-11.99	0.00	0.00

(d) QTRAN

$a_0 \backslash a_1$	A	B	C
A	8.00	-12.00	-12.00
B	-11.99	0.00	0.00
C	-11.99	0.00	0.00

(e) RQN

#### D. Training the Residual Q-Network

Similar to the description of DQN as in (1), the networks that we use during learning are trained to minimize the TD loss that is now generalized for  $Q_{tot}^+$  as

$$\mathcal{L}_b(\theta_b) = \mathbb{E} \left[ (y_b^{tot} - Q_{tot}^+(\tau, a; \theta))^2 \right], \quad b \in \{1, \dots, B\} \quad (13)$$

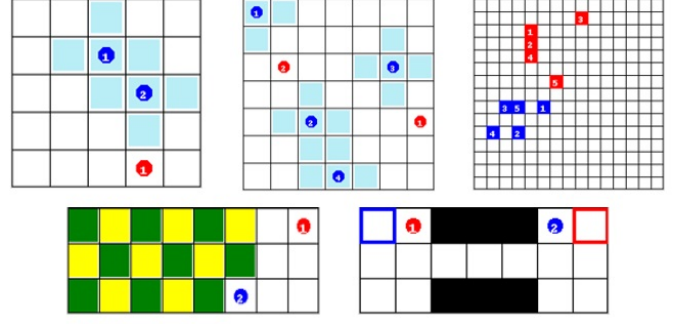


Fig. 3: Representation of the matrix environments used in the experiments [30]. On the top, from the left to the right: PredatorPrey\_2, PredatorPrey\_4, Combat. On the bottom, from the left to the right: Checkers, Switch.



Fig. 4: Representation of the Starcraft II environments used in the experiments. From the left to the right: 8m, 3s5z, 3s5z\_vs\_3s6z.

where  $B$  is the number of batches sampled from the replay buffer and  $y_b^{tot} = R + \gamma \max_{a'} Q_{tot}(\tau', a'; \theta^-)$  is the target for an iteration  $b$ , where  $\theta^-$  represents the parameters of the target network. This loss is propagated firstly through the entire evaluation network and then through the RQN. As such, the targets and the evaluation values are calculated in a slightly different way. The target network is based on  $Q_{tot}$  instead of  $Q_{tot}^+$ . In other words, the target network of the proposed architecture is based on the output of a VDN network,  $Q_{tot}$ , as described previously in subsection III-B. The motivation for this is that by keeping a separate set of targets, we can approximate them with successively minimized losses due to our estimation factor. The RQN will eventually reach a state where it is no more needed since the Q-values will be optimal and stable. At that point, it will render itself obsolete. (this is a key difference of the proposed RQN from ResNets [28]). Thus, the RQN acts as an auxiliary acceleration network that will encourage the agents to converge quickly, by aiding the decomposition of the joint action-value function.

## V. EXPERIMENTAL DETAILS<sup>1</sup>

### A. Details on the Environments Used

The performance of the proposed method is evaluated in four different matrix environments with distinct objectives and Starcraft II (SC2) [31] (Figure 4). The maps of the first set of environments are 2D matrices (Figure 3) where the action space is discrete, and the observation space is continuous. This set of environments was chosen since similar versions of them were used in previous works such as QTRAN [14] and VDN [13]. Note that in the described environments, the individual

<sup>1</sup>Code at <https://github.com/rafaelmp2/residual-q-net>

rewards attributed are then summed up to form a shared team reward.

**Switch:** two agents start in opposite cells of the map, separated by a narrow corridor and must switch positions. This environment helps to evaluate how the agents react to a delayed reward scheme. **Checkers:** on a map with apples and lemons, the agents must eat all the existing apples on the map and avoid lemons. For that, the agents need to cooperate to maximize the team reward. **PredatorPrey:** a group of predators must catch a group of moving prey. There is a step penalty of -0.01 and every time one of the prey is caught, each agent in the team receives a reward of +5. At least two agents are needed to catch one prey. Two different versions of this environment are used: one where there are 2 predators and 1 prey and in the other, there are 4 predators and 2 prey. These are referred to as PredatorPrey\_2 and PredatorPrey\_4, respectively. Since this environment only has a success reward and no significant intermediate penalties, the task remains monotonic [14]. **Combat:** a team of five agents must learn to cooperate to eliminate all the agents of the enemy team. **SC2:** to complement our experiments we used environments from the challenging Starcraft Multi-Agent Challenge (SMAC) set of environments [31]. The maps used are the *8m*, *3s5z*, *3s5z\_vs\_3s6z*, all with 8 agents. The level of difficulty was set to very difficult (level 7).

### B. Algorithms Used as Benchmarks

The performance of the proposed method is benchmarked against other related approaches using the average episode reward as metric. The agents are trained for a different number of episodes depending on the complexity of the environment, each in 3 independent runs starting at different random seeds. During training time, at every 100 episodes, the model performance is evaluated by running 20 episodes and by averaging the joint rewards for each episode. The proposed approach is compared against VDN [13], QMIX [12], QTRAN-base [14], QPLEX [17] and Weighted QMIX (WQMIX) [18]. The width of the hidden layers is 32 for QMIX and 64 for RQN and the other methods, as described in the respective papers. The rest of the hyperparameters used in the main experiments are summarized in II.

## VI. RESULTS AND DISCUSSIONS

In this section we present now the results of the experiments described in the previous section. We start by discussing the results in the matrix environments and analyse the stability of the learning process. Then we discuss the results in the Starcraft II environments and finally we present additional results to analyse the algorithm performance.

### A. Performance in Matrix Environments

Figure 5 illustrates the evolution of the performance of the proposed Residual Q-Network (abbreviated as RQN) and the benchmarked methods. The curves represent the smoothed average rewards over the episodes (smoothed with a 10-step cumulative moving average). For the Switch environment

(Figure 5e), RQN and WQMIX methods manage to solve the task in a small amount of time. Considering that this is a task with delayed rewards, this enforces the intuition that the proposed RQN method encourages the agents to prefer promising states once they are explored. In other words, since the agents do not receive intermediate information about whether their actions are good or not, once they get to the goal, the proposed method encourages them to visit that state again successively, which explains the stability once an optimal state is achieved, making RQN method outperform VDN over time. In the additional experiments subsection it will be clearly illustrated the concept of stability. VDN and QPLEX also solves this task, but are unable to do it in the same amount of time as the previous, while QMIX and QTRAN-base fail to solve it.

The results in 5a show that the RQN method can solve a complex non-monotonic task like Checkers successfully. Furthermore, we can see that it is capable of maintaining a stable performance over time once it learns to solve the task. In comparison, although QMIX, VDN, QPLEX QTRAN-base and WQMIX manage to solve the task, they fail to achieve as high rewards and stable performance as RQN.

Analysing the performances obtained for two versions of PredatorPrey, all the proposed methods can achieve optimal rewards. In PredatorPrey\_2 (Figure 5b), all the evaluated methods are able to maintain a stable optimal reward over time. Moreover, since this is a monotonic task, QMIX can solve it successfully as it was expected. However, when the number of agents is increased and hence the complexity of the environment, in PredatorPrey\_4 (Figure 5c) we observe a different scenario. Despite all the methods being capable of achieving an high value, VDN shows an unexpected behaviour in this task with a notable decrease of performance after around 5000 episodes of training. We believe this happens because the agents get confused due to the exploration activities of the other agents. With the proposed method we aim to solve this problem by improving the credit assignment of the agents and consequently reducing the relative overgeneralization. The results suggest that this problem does not incur in RQN. Surprisingly, QTRAN-base and WQMIX stay slightly below the other methods, while QPLEX also solves this task successfully achieving high rewards. However, RQN can still maintain a higher and more consistent performance. Combat is the exception to the described pattern in the previous environments, where RQN stays below the performance of VDN or WQMIX in the final training steps. On the other hand, QTRAN-base fails to solve this environment and as expected, since this task is not monotonic, QMIX also fails to solve it.

### B. Stability of Performance in the Learning Process

In Figure 5 we can identify a pattern in all the environments (with the slight exception of Combat) where RQN shows a stable result in the reward over time, opposing to the other methods. Although most of them are able to solve the tasks, they fail to achieve as high values as RQN and show a lack of stability over time, due to high variance of the rewards obtained and to failing to consistently achieve an optimal

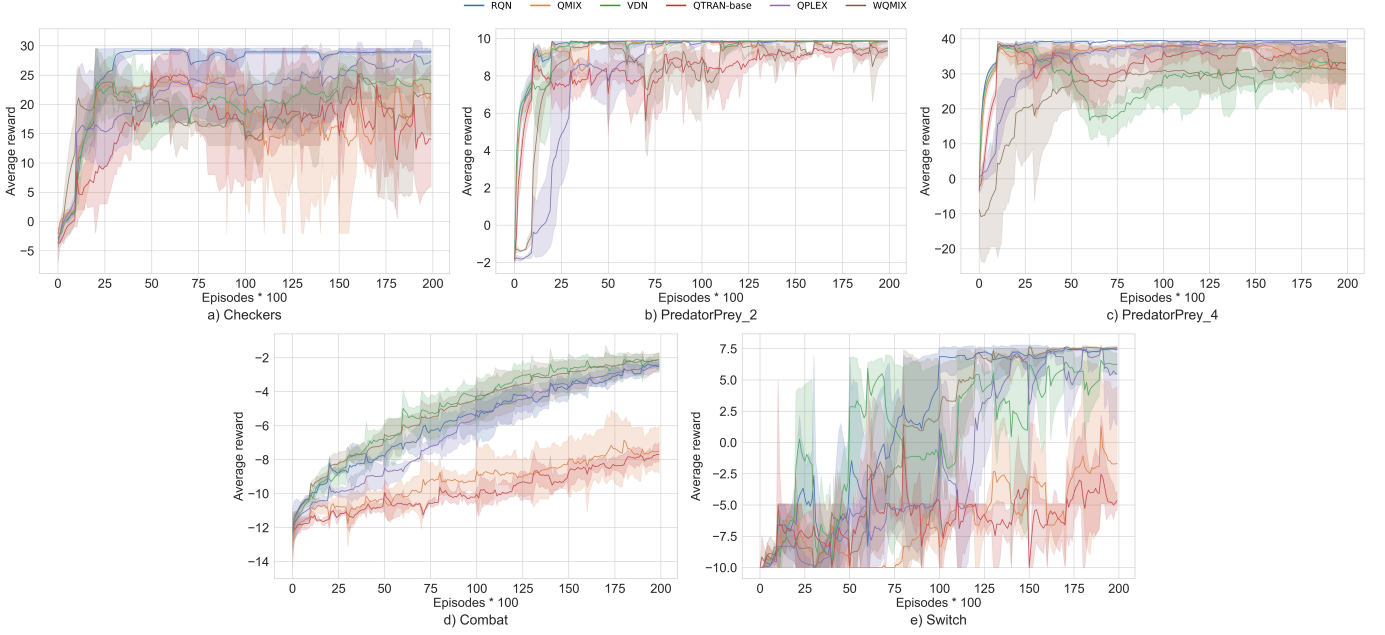


Fig. 5: Smoothed average rewards for the experimented methods in the matrix environments. The bold area represents the 95% confidence intervals.

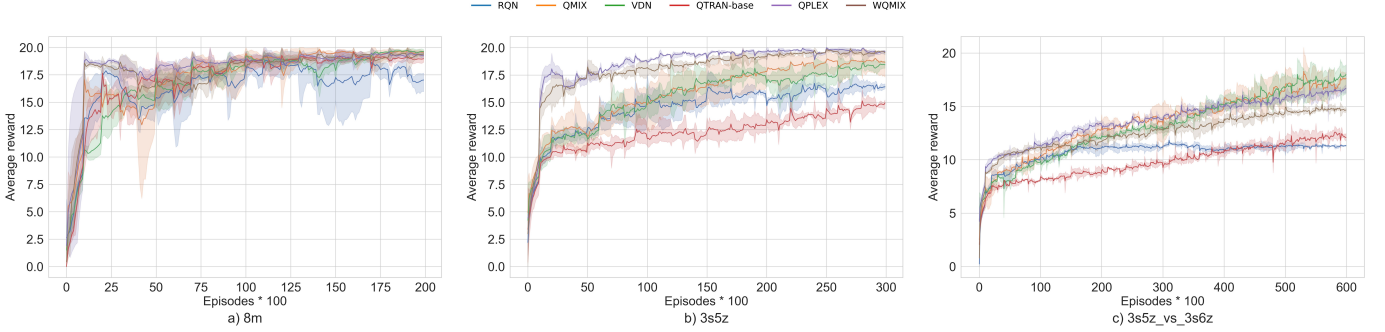


Fig. 6: Smoothed average rewards in the attempted Starcraft II environments. The bold area represents the 95% confidence intervals.

reward. While RQN is able to maintain the performance achieving high rewards over time, we can see that the other methods keep decaying. In Switch, RQN is also the only one able to maintain a stable performance over time. We call this effect the stability of the learning process. We believe the exceptions with Combat are due to the complexity of the environment allied to its very delayed rewards. Generally, in Figure 5 it is possible to see that RQN displays a very reduced number of oscillations when the convergence (learning goal) was achieved. This is a key difference from the other methods that fail to maintain a stable performance, even when the learning objective is achieved. Observing now PredatorPrey\_4 (Figure 5c), we can see that QMIX achieves high performance but shows it to be very unstable over time in this four-agent task. This observation suggests that, although the task being monotonic, QMIX fails to scale with the increase in the number of agents. In comparison to QPLEX, although it is capable of maintaining optimal and stable values over time, the proposed RQN method can achieve the same optimal values

and, at the same time, shows improved level of stability.

### C. Performance in Starcraft II Environments

Figure 6 illustrates the average rewards in the SC2 environments (smoothed with a 10-step cumulative moving average). The corresponding win rates can be found in the Appendix A. Although RQN may not outperform the other methods in this set of environments, it shows to be competitive. In the attempted maps, like the majority of the compared methods, RQN can learn the 3s5z and 8m tasks and also show competitive performance in the 3s5z\_vs\_3s6z, although none of them is able to solve the last. Comparing to QTRAN-base RQN can still stay above in the last two tasks. These results, linked to the previous in figures 5, support the fact that RQN can perform in a wide family of environments due to its relaxed factorization.

Overall, we can see that QTRAN could not learn a good factorization for Combat and Switch environments. QMIX and VDN also fail when monotonicity and additivity are not met,



TABLE II: Hyperparameters used in the experiments in section VI.

HYPERPARAMETER	VALUE	DESCRIPTION
BATCH SIZE	32	NUMBER OF SAMPLES AT EACH TRAINING STEP
BUFFER SIZE	5000	MAXIMUM NUMBER OF SAMPLES TO STORE IN THE BUFFER
INITIAL $\epsilon$	1	STARTING VALUE FOR $\epsilon$
MINIMUM $\epsilon$	0.05	FINAL VALUE FOR $\epsilon$
$\epsilon$ ANNEAL STEPS	50000	NUMBER OF STEPS THAT $\epsilon$ TAKES TO DECAY

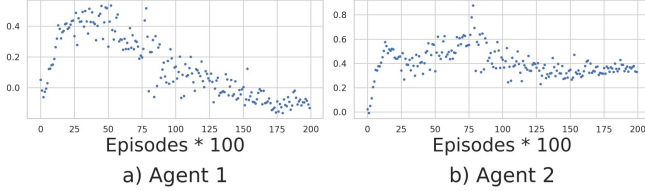


Fig. 7: Estimation factors learned by the RQN over time for the two agents in PredatorPrey\_2.

respectively [14]. When compared to QPLEX and WQMIX, the results suggest that RQN is capable of learning tasks with high performance in a more variate number of environments. More specifically, RQN proves to be competitive in the SC2 environments attempted, linked to a superior performance in the tasks illustrated in Figure 5. This suggests that RQN can learn a more variate number of different environments, and specially in the absence of extra state information during training. Methods that use a relaxed factorization like RQN and QTRAN might not perform as well in environments such as SC2 [17]. However, a good relaxation proves to be very effective in environments where non-cooperative behaviours are punished more severely [14], such as the ones in Figure 5. The results in Figures 5 and 6 confirm that, with a more relaxed factorization, RQN can factorize a wider family of environments, even when compared to QTRAN that also uses a relaxed and unconstrained factorization.

#### D. RQN Self-Depreciation and Task Monotonicity

In this subsection, additional experiments on the algorithm performance are presented to complement the main results in section VI.

1) *Estimation Factors Learned by the RQN*: Figure 7, illustrates the estimation factors learned by the RQN over time in the PredatorPrey\_2 environment. It is possible to see that it learns an estimation factor over time, and when the performance objectives are reached, the estimation factor for each agent converges to a stable value. Once it reaches the stable value, the purpose of RQN as an auxiliary network that assists in value factorization becomes obsolete. We call this effect the self-depreciation of the RQN when the training objective is reached.

Figure 8 shows the same metric for the PredatorPrey\_4 environment (with 4 agents). This environment is more complex due to the increased number of agents. Thus, although we can see the estimation factor clearly converging for agents 1 and

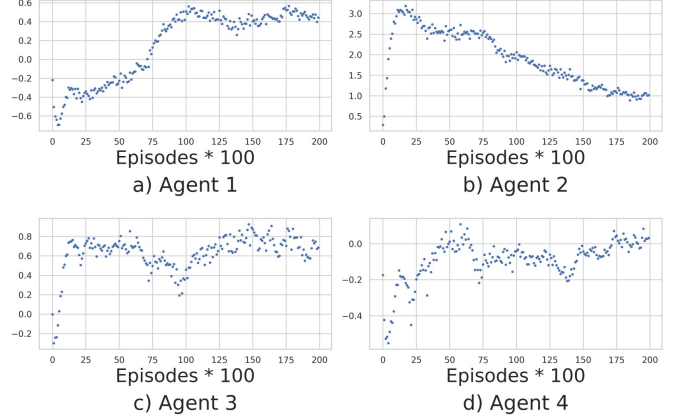


Fig. 8: Estimation factors learned by the RQN over time for the four agents in PredatorPrey\_4.

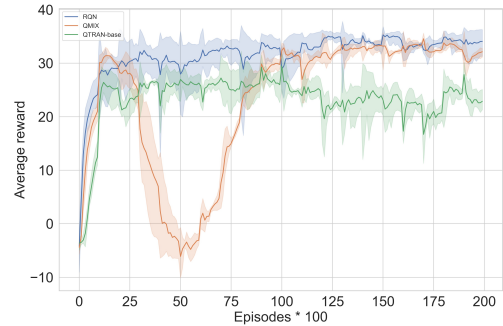


Fig. 9: Smoothed average rewards for PredatorPrey\_4. The bold area represents the 95% confidence intervals.

2, with agent 3 and 4 the convergence has not stabilized. This suggests that, in this case, the agents are still exploring. These representations show how the RQN self-deprecates, because when the training objective is reached, the values estimated by this network should remain stable. The reason is that at the moment of convergence and the objective was achieved, the individual trajectories no longer need to be adjusted.

2) *Impact of a Lower Level of Monotonicity*: To evaluate the performance of QTRAN, QMIX and RQN against a change in the level of the monotonicity in the tasks experimented in this paper, we show now the results for the same PredatorPrey task with 4 agents as in VI but now with a capture penalty of -0.1. This means that, if one of the agents tries to capture the prey alone, then each agent of the team receives a penalty of -0.1 (each agent gets this penalty that will be added into a

TABLE III: Hyperparameters used in subsubsection VI-D2 for the modified PredatorPrey\_4.

HYPERPARAMETER	VALUE	DESCRIPTION
BATCH SIZE	32	NUMBER OF SAMPLES AT EACH TRAINING STEP
BUFFER SIZE	100000	MAXIMUM NUMBER OF SAMPLES TO STORE IN THE BUFFER
INITIAL $\epsilon$	1	STARTING VALUE FOR $\epsilon$
MINIMUM $\epsilon$	0.1	FINAL VALUE FOR $\epsilon$
$\epsilon$ ANNEAL STEPS	2000000	NUMBER OF STEPS THAT $\epsilon$ TAKES TO DECAY

shared team reward). The results are smoothed with a 10-step cumulative moving average (Figure 9).

In Figure 9 we can see that RQN can still maintain a good performance over time, while the benchmarked methods suffer from this change on the environment, mainly QMIX. The results suggest that QMIX is more sensitive to monotonicity changes in the environment as it was proved before in [14]. RQN however is capable of still performing well even with a change in the level of monotonicity of the environment, suggesting that RQN eliminates this constraint imposed by QMIX. This observation is also supported by some of the main experiments in section VI with some non-monotonic tasks. The hyperparameters used in this additional experimental setting are summarized in Table III.

## VII. CONCLUSIONS AND FUTURE WORK

We study the problem of value-function factorization for MARL within the scope of centralized training and decentralized execution scheme [8], [9]. Here, we propose a method that uses a novel agent configuration, named as Residual Q-Networks (RQNs), which transforms the individual factor functions in a way that preserves the Individual-Global-Max (IGM) property, and at the same time enables the factorization in a robust way. The RQN is designed as an auxiliary network that learns to adjust the action-value function trajectories of individual agents to reflect the relative importance of different episodes. With the proposed RQN based approach, it is possible to factorize a wider family of multi-agent environments by improving the credit assignment problem [7], [16], an issue that is prominent in decentralized learning. The proposed RQN architecture not only achieves better performance, but also improves the performance stability over time.

The experimental results illustrate that the proposed method can outperform the state-of-the-art methods such as QTRAN, VDN, QMIX, QPLEX and WQMIX in certain tasks with strong punishments for non-cooperative behaviours and with varying levels of complexity and with. An exception to this is in StarCraft environments, where complete state information is available, but RQN is yet competitive and learns to solve the task. The results suggest that the proposed RQN based approach is capable of factorizing value-functions for more families of environments than the previous method, owing to a more robust and relaxed transformation that preserves IGM, and a neural architecture that is more tractable. Moreover, opposing to most of the benchmarked methods, RQN does not take advantage of full state information during the training phase. This enhances the practical application of RQN, especially towards non-simulated environments where complete

state information is unavailable even during training time. Future work involves the scalability of MARL towards more complex environments, where the number of agents gradually increases. Accordingly, we aim to investigate how curriculum-based approaches can be used to extend the current method.

## APPENDIX WIN RATES FOR THE STARCRAFT II ENVIRONMENTS

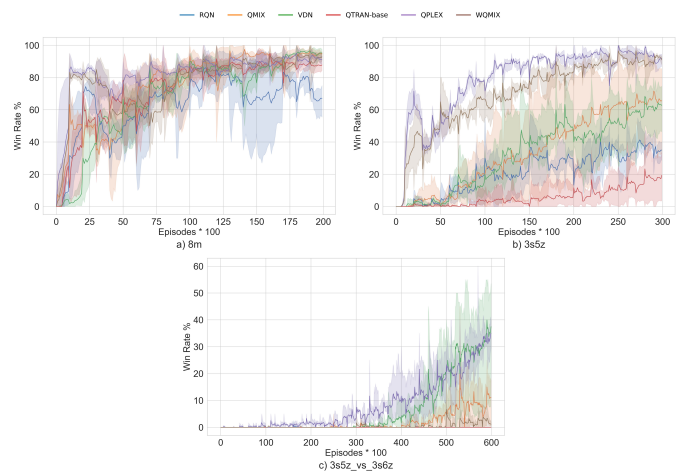


Fig. 10: Win rates for the experimented methods in the Starcraft II environments. The bold area represents the 95% confidence intervals.

## REFERENCES

- [1] G. Chen, “A New Framework for Multi-Agent Reinforcement Learning – Centralized Training and Exploration with Decentralized Execution via Policy Distillation,” in *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, 2020, p. 1801–1803.
- [2] C. Qu, S. Mannor, H. Xu, Y. Qi, L. Song, and J. Xiong, “Value Propagation for Decentralized Networked Deep Multi-Agent Reinforcement Learning,” in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, 2019, pp. 1182–1191.
- [3] J. Hook, S. El-Sedky, V. De Silva, and A. Kondoz, “Learning data-driven decision-making policies in multi-agent environments for autonomous systems,” *Cognitive Systems Research*, vol. 65, pp. 40–49, 2021.
- [4] S. Iqbal and F. Sha, “Actor-Attention-Critic for Multi-Agent Reinforcement Learning,” in *Proceedings of the 36th International Conference on Machine Learning*, vol. 97. PMLR, 09–15 Jun 2019, pp. 2961–2970.
- [5] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, “Deep Reinforcement Learning for Multi-Agent Systems: A Review of Challenges, Solutions and Applications,” *IEEE Transactions on Cybernetics*, vol. 50, no. 9, pp. 3826–3839, Sep. 2020.
- [6] Z. Zhang, D. Wang, and J. Gao, “Learning Automata-Based Multiagent Reinforcement Learning for Optimization of Cooperative Tasks,” *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–14, 2020.

- [7] E. Wei, D. Wicke, D. Freelan, and S. Luke, "Multiagent Soft Q-Learning," *AAAI Spring Symposium Series*, 2018.
- [8] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, p. 6382–6393.
- [9] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual Multi-Agent Policy Gradients," in *AAAI 2018: Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, February 2018.
- [10] J. K. Gupta, M. Egorov, and M. Kochenderfer, "Cooperative Multi-Agent Control Using Deep Reinforcement Learning," in *Autonomous Agents and Multi-Agent Systems*, 2017, vol. 10642, pp. 66–83.
- [11] J. Hook, V. D. Silva, and A. Kondoz, "Deep Multi-Critic Network for accelerating Policy Learning in multi-agent environments," *Neural Networks*, vol. 128, pp. 97–106, 2020.
- [12] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson, "QMIX: Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning," in *Proceedings of the 35th International Conference on Machine Learning*, vol. 80, 10–15 Jul 2018, pp. 4295–4304.
- [13] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. F. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, and T. Graepel, "Value-Decomposition Networks for Cooperative Multi-Agent Learning Based on Team Reward," in *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS 2018, Stockholm, Sweden, July 10-15, 2018*, 2018, pp. 2085–2087.
- [14] K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi, "QTRAN: Learning to factorize with transformation for cooperative multi-agent reinforcement learning," in *Proceedings of the 36th International Conference on Machine Learning*, vol. 97, 09–15 Jun 2019, pp. 5887–5896.
- [15] B. J. Lansdell, P. R. Prakash, and K. P. Kording, "Learning to solve the credit assignment problem," in *International Conference on Learning Representations*, 2020.
- [16] J. Yang, A. Nakhaei, D. Isele, K. Fujimura, and H. Zha, "CM3: Cooperative Multi-goal Multi-stage Multi-agent Reinforcement Learning," in *International Conference on Learning Representations*, 2020.
- [17] J. Wang, Z. Ren, T. Liu, Y. Yu, and C. Zhang, "QPLEX: Duplex Dueling Multi-Agent Q-Learning," in *International Conference on Learning Representations*, 2021.
- [18] T. Rashid, G. Farquhar, B. Peng, and S. Whiteson, "Weighted QMIX: Expanding Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 10 199–10 210.
- [19] M. Tan, "Multi-Agent Reinforcement Learning: Independent vs. Cooperative Agents," in *In Proceedings of the Tenth International Conference on Machine Learning*, 1993, pp. 330–337.
- [20] C. Watkins and P. Dayan, "Technical Note: Q-Learning," *Machine Learning*, vol. 8, pp. 279–292, 1992.
- [21] A. Tampuu, T. Matiisen, D. Kodelja, I. Kuzovkin, K. Korjus, J. Aru, J. Aru, and R. Vicente, "Multiagent cooperation and competition with deep reinforcement learning," *PLOS ONE*, vol. 12, no. 4, p. e0172395, 2017.
- [22] C. Sun, W. Liu, and L. Dong, "Reinforcement learning with task decomposition for cooperative multiagent systems," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 5, pp. 2054–2065, 2021.
- [23] J. Foerster, I. A. Assael, N. de Freitas, and S. Whiteson, "Learning to Communicate with Deep Multi-Agent Reinforcement Learning," in *Advances in Neural Information Processing Systems*, vol. 29, 2016.
- [24] S. Sukhbaatar, A. Szlam, and R. Fergus, "Learning Multiagent Communication with Backpropagation," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016, p. 2252–2260.
- [25] A. Das, T. Gervet, J. Romoff, D. Batra, D. Parikh, M. Rabbat, and J. Pineau, "TarMAC: Targeted Multi-Agent Communication," in *Proceedings of the 36th International Conference on Machine Learning*, vol. 97, 09–15 Jun 2019, pp. 1538–1546.
- [26] F. A. Oliehoek and C. Amato, *A Concise Introduction to Decentralized POMDPs*, 1st ed. Springer Publishing Company, Incorporated, 2016.
- [27] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529–533, 02 2015.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *2016 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- [29] K. He, X. Zhang, S. Ren, and J. Sun, "Identity Mappings in Deep Residual Networks," in *ECCV (4)*, 2016, pp. 630–645.
- [30] A. Koul, "ma-gym: Collection of Multi-Agent environments based on OpenAI Gym," *GitHub repository*, 2019. [Online]. Available: <https://github.com/koulaurag/ma-gym>
- [31] M. Samvelyan, T. Rashid, C. Schroeder de Witt, G. Farquhar, N. Nardelli, T. G. J. Rudner, C.-M. Hung, P. H. S. Torr, J. Foerster, and S. Whiteson, "The StarCraft Multi-Agent Challenge," in *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, 2019.

**Rafael Pina** is currently a Ph.D. student at Loughborough University London in machine intelligence. His work focus mainly on Multi-Agent Reinforcement Learning. He got his B.S. in Informatics Engineering at the University of Coimbra in 2019 and his M.S. in Cyber Security and Big Data at Loughborough University London in 2020.

**Varuna De Silva** received his Ph.D. from University of Surrey in United Kingdom in 2011. He was a postdoctoral researcher in the same institute between 2011-2013, working on video and image processing for Multiview video broadcasting applications. Since 2013 November, he worked as a senior algorithms developer for image signal processors at Apical Ltd, UK (Now part of ARM Plc). In April 2016, he joined Loughborough University, as a lecturer in digital engineering, where he was promoted to a Senior Lecturer in Jan 2020.

**Joosep Hook** holds a B.S. in Computer Engineering from the University of Tartu. He is currently pursuing his Ph.D. in machine learning at Loughborough University London. His work focus on Multi-Agent Reinforcement Learning and Curriculum Learning.

**Ahmet Kondoz** received the B.S. (Hons.) degree in engineering from the University of Greenwich, Greenwich, U.K., in 1983, and the Ph.D. degree in telecommunication from the University of Surrey, Guildford, U.K., in 1987. He became a Lecturer in 1988, a Reader in 1995, and then a Professor in Multimedia Communication Systems in 1996, at the University of Surrey. Ahmet took part in setting up of the world renowned Centres CCSR and the I-Lab at the University of Surrey before joining Loughborough University London where he is now the Director of the Institute for Digital Technologies. He has published more than 400 journal and conference papers, three books, and nine patents.