# Locality Matters: A Scalable Value Decomposition Approach for Cooperative Multi-Agent Reinforcement Learning

**Roy Zohar[1], Shie Mannor[2], Guy Tennenholtz[2]**

[1] Hebrew University of Jerusalem
[2] Nvidia Research, Technion Institute of Technology
roy.zohar@mail.huji.ac.il, shie@ee.technion.ac.il, guytenn@gmail.com

## Abstract

Cooperative multi-agent reinforcement learning (MARL) faces significant scalability issues due to state and action spaces that are exponentially large in the number of agents. As environments grow in size, effective credit assignment becomes increasingly harder and often results in infeasible learning times. Still, in many real-world settings, there exist simplified underlying dynamics that can be leveraged for more scalable solutions. In this work, we exploit such locality structures effectively whilst maintaining global cooperation. We propose a novel, value-based multi-agent algorithm called LOMAQ, which incorporates local rewards in the Centralized Training Decentralized Execution paradigm. Additionally, we provide a direct reward decomposition method for finding these local rewards when only a global signal is provided. We test our method empirically, showing it scales well compared to other methods, significantly improving performance and convergence speed.

## 1 Introduction

The field of Reinforcement Learning (RL) is concerned with an agent taking actions in an environment in order to maximize a cumulative reward. Recent work has witnessed major success in various tasks, including Atari games (Mnih et al. 2015), and Go (Silver et al. 2016). A popular extension of RL is cooperative multi-agent RL (cooperative MARL), in which a group of agents attempts to interact with an environment together. Research on MARL has gained much attention in recent years, with examples in the Star-Craft multi-agent challenge (Vinyals et al. 2019) and traffic control (Chu et al. 2019).

A common paradigm used in cooperative MARL is Centralized Training Decentralised Execution (CTDE, Kraemer and Banerjee (2016)). In this approach, agents are trained simultaneously by a centralized controller. Decentralized policies are then derived from the training process and used for execution. Centralized training can be highly beneficial, granting access to additional global information, which helps agents coordinate their actions. Nevertheless, utilizing such information effectively is a challenging problem for cooperative MARL, due to exponential state and action
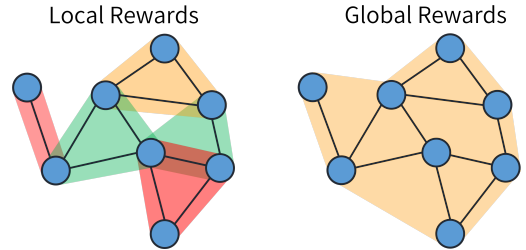


Figure 1: A visualization of training MARL for a graph of agents with local rewards vs. with global rewards. The colored regions represent the feedback that the agents exhibit during training

spaces. As the environment scales, coordination becomes increasingly difficult, rendering centralized training impractical. Still, in many real-world settings, there exist simplified underlying dynamics that can help tackle this problem.

In this paper, we utilize *local rewards*, a principal component of our work. While local rewards are often used in competitive settings (i.e., where every agent attempts to maximize its own local reward), in most cooperative approaches, cooperation is weakly enforced through a shared global reward that all agents aim to maximize (Rashid et al. 2018; Lowe et al. 2017). A visualization of this paradigm is depicted in Figure 1

Local rewards are critical for effective learning in scalable settings. As an example, consider the problem of coaching a large soccer team. If a certain player loses the ball to the other team, punishing that player directly (and possibly neighboring players) with targeted feedback, may be far more effective than punishing the entire team with general feedback. The latter will often leave players confused, believing they should have acted differently.

Despite the effectiveness of local rewards, naively training with local rewards may result in greedy agents that fail to cooperate. Concurrent approaches that aim to exploit local reward structures for our setting often pay a price in terms of cooperation and usually resort to training with global rewards (Lowe et al. 2017). This is particularly true for the value decomposition approach for cooperative MARL, which has become increasingly popular in recent

years (Sunehag et al. 2018; Rashid et al. 2018; Son et al. 2019; Rashid et al. 2020; Wang et al. 2020). To the best of our knowledge, there are no value decomposition methods that utilize local rewards effectively. Rather, they rely on the global reward signal for decomposing the joint state-action value function into individual state-action value functions. As we show in our work, such an approach hurts overall performance and convergence speed in large environments.

In this work, we present a scalable value decomposition method for the cooperative CTDE setting. Our method leverages local agent rewards for improving credit assignment, whilst maintaining a cooperative objective. In addition, we provide a direct decomposition method for finding local rewards when only a global reward is provided. We empirically show that our method is scalable, improving upon state-of-the-art methods for this setting.

Our contributions are as follows. We define the $Q$-Summation Maximization (QSM) Condition (Section 3.1), showing its theoretical benefits in a linear bandit setting (Theorem 1). We show that a monotonic decomposition of utilities can be derived to establish the QSM condition (Section 3.3), and provide a value-based algorithm to enforce it (Section 4). Finally, we construct a reward decomposition method for learning local rewards when a global reward is given (Section 4.2).

## 2 Preliminaries

We define a multi-agent Markov decision process (MAMDP) as the tuple $\mathcal{M} = (\mathcal{G}, \mathcal{S}, \mathcal{A}, P, r, \gamma)$, where $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is an undirected graph of agents, where $\mathcal{V} = [n] = \{1, \ldots, n\}$ and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$, $\mathcal{S} = \times_{i=1}^n \mathcal{S}_i$ is the global state space, $\mathcal{A} = \times_{i=1}^n \mathcal{A}_i$ is the global action space, $P : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ is the global transition function, $r : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ is the global reward, and $\gamma \in (0, 1)$ is the discount factor.

An agent $i \in \mathcal{V}$ is associated with the underlying graph $\mathcal{G}$, state $s_i$ and action $a_i$. For a set $B \subseteq \mathcal{V}$ we define $s_B, a_B$ as the subset of agent states and actions in $B$, i.e., $s_B = (s_i)_{i \in B}$ and $a_B = (a_i)_{i \in B}$, respectively. At time $t$, the environment is at state $s = (s_1, \ldots, s_n)$ and the agents take an action $a = (a_1, \ldots, a_n)$, after which the environment returns a reward $r$ and transitions to state $s'$ according to the factored dynamics $P(s'|s, a) = \prod_{i \in \mathcal{V}} P_i(s'_i|s_{N(i)}, a_i)$, where here we used $N(i)$ to denote the neighborhood of agent $i$, including $i$, i.e., $N(i) = \{j \in \mathcal{V} : (i, j) \in \mathcal{E}\} \cup \{i\}$.

We define a global Markovian policy $\pi$ as a mapping $\pi : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ such that $\pi(a|s)$ is the probability to choose action $a = (a_0, \ldots, a_n)$ at state $s = (s_0, \ldots s_n)$. We define the value of policy $\pi$ starting at a state $s \in \mathcal{S}$ and taking action $a \in \mathcal{A}$ as

$$Q^\pi(s, a) = \mathbb{E}_\pi \left[ \sum_{t=0}^\infty \gamma^t r(s(t), a(t)) \,\middle|\, s(0) = s, a(0) = a \right].$$

The value function is then defined by $v^\pi(s) = \mathbb{E}_{a \sim \pi(s)}[Q^\pi(s, a)]$. We define the optimal value and optimal policy by $v^*(s) = \max_\pi v^\pi(s)$ and $\pi^* \in \arg\max_\pi v^\pi(s)$, respectively.

Finally, we denote by $\mathcal{P}$ a partition of $\mathcal{V} = [n]$ (i.e., of agents), such that $\bigcup_{J \in \mathcal{P}} J = \mathcal{V}$ and $\bigcap_{J \in \mathcal{P}} J = \emptyset$. We say that $\mathcal{P}'$ is a refinement of $\mathcal{P}$ if for every $J' \in \mathcal{P}'$ there exists $J \in \mathcal{P}$ such that $J' \subseteq J$[1].

### 2.1 Reward Decomposition

A primary element of MARL is the decomposition of the reward function $r$ over agent states and actions $\{s_i, a_i\}_{i \in \mathcal{V}}$. Given some decomposition of rewards $\{r_i : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}\}_{i \in \mathcal{V}}$, such that $r(s, a) = \sum_{i \in \mathcal{V}} r_i(s, a)$, we define the partial $Q$-function of $\pi$, denoted by $Q_i^\pi(s, a)$ as $Q_i^\pi(s, a) = \mathbb{E}_\pi[\sum_{t=0}^\infty \gamma^t r_i(s(t), a(t)) \mid s(0) = s, a(0) = a]$. It follows that $Q^\pi(s, a) = \sum_{i \in \mathcal{V}} Q_i^\pi(s, a)$. Note that such decomposition always exists, e.g., by choosing $r_1 = r, r_i = 0, i \geq 2$.

In this work, we consider a reward decomposition for which every agent is dependent only on its local state and action, as defined formally below. We refer the reader to Section 4.3 for a relaxation of this assumption.

**Assumption 1** (Qu et al. (2020)). *We assume that the reward function $r$ is additively decomposable. That is, there exist $\{r_i : \mathcal{S}_i \times \mathcal{A}_i \mapsto \mathbb{R}\}_{i \in \mathcal{V}}$ such that $r(s, a) = \sum_{i=1}^n r_i(s_i, a_i)$ for all $s = (s_1, \ldots, s_n)$, $a = (a_1, \ldots, a_n)$.*

## 3 Value Partitions for MARL

In this section, we focus on leveraging value-based partitions for credit assignment in MARL. We consider decoupling the problem into smaller problems, each of which can be viewed as a separate, easier estimation problem. Particularly, we generalize ideas from Rashid et al. (2018), and define a partition-based $Q$-maximization condition. We motivate this condition in a contextual bandit setting, proving it can exponentially improve regret. Then, for the general RL setting, we propose a monotonic decomposition of agent utilities for which our proposed condition holds. We show examples of the latter and prove that monotonic decomposition of utilities is indeed sufficient for partition-based maximization. Our decomposition will prove beneficial in Section 4, as we leverage value partitions to construct a scalable value-based algorithm for MARL.

### 3.1 $Q$-Summation Maximization (QSM)

We begin by defining the $Q$-Summation Maximization condition on which we build upon the rest of this section. The QSM condition states that the $Q$-function can be maximized using a partition of partial maximizers, as defined formally below.

**Definition 1** (QSM Condition). *Let $\mathcal{P}$ be a partition of $\mathcal{V}$. We say that a MAMDP satisfies the Q-Summation Maximisation (QSM) Condition with $\mathcal{P}$, if for every $s \in \mathcal{S}$ and policy $\pi$*

$$\max_a \left\{ \sum_{i=1}^n Q_i^\pi(s, a) \right\} = \sum_{J \in \mathcal{P}} \left( \max_a \left\{ \sum_{i \in J} Q_i^\pi(s, a) \right\} \right)$$

[1]A refinement partition can be useful when multiple groups of agents concurrently attempt to solve relatively separable tasks.

**Algorithm 1: Multi-OFUL**

1: **input:** $\alpha, \lambda, \delta > 0$, $\mathcal{P}$ partition of $\mathcal{V}$
2: **init:** $V_{J,a_J} = \lambda I$, $J \in \mathcal{P}$, $a_J \in \times_{i \in J} \mathcal{A}_i$.
3:    $Y_J = 0$, $J \in \mathcal{P}$.
4: **for** $t = 1, 2, \ldots$ **do**
5:    Receive context $x(t)$
6:    **for** $J \in \mathcal{P}$, $a_J \in \times_{i \in J} \mathcal{A}_i$ **do**
7:       $\hat{y}_{a_J}(t) = \left\langle x(t), V_{J,a_J}^{-1} Y_J \right\rangle$
8:       $\text{UCB}_{a_J}(t) = \sqrt{\beta_J(t, \delta)} \|x(t)\|_{V_{J,a_J}^{-1}}$
9:    **end for**
10:    $a(t) \in \times_{a_J} \arg\max_{J \in \mathcal{P}} \hat{y}_{a_J}(t) + \alpha \text{UCB}_{a_J}(t)$
11:    Play $a(t)$ and observe $\{r_J(t)\}$
12:    $V_{J,a_J(t)} = V_{J,a_J(t)} + x(t)x(t)^T$, $J \in \mathcal{P}$
13:    $Y_J = Y_J + x(t) r_J(t)$, $J \in \mathcal{P}$
14: **end for**

We note the two extremes of the QSM Condition. First, every MAMDP satisfies the condition trivially with $\mathcal{P} = \{\mathcal{V}\}$. Second, if $\mathcal{P}$ partitions $\mathcal{V}$ into singletons (i.e., $\mathcal{P} = \{\{1\}, \{2\}, \ldots, \{n\}\}$), then for every $s \in \mathcal{S}$,

$$\max_a \left\{ \sum_{i=1}^n Q_i^\pi(s, a) \right\} = \sum_{i=1}^n \left( \max_a \{Q_i^\pi(s, a)\} \right).$$

The QSM condition can greatly improve learning efficiency in settings in which the partial $Q$-functions are easier to approximate, effectively decoupling the problem to $|\mathcal{P}|$ simpler problems. We prove this for an instance of the linear bandits problem in the following subsection. Then, in Section 3.3 we discuss a sufficient assumption for which the QSM condition holds.

### 3.2 QSM in Linear Bandits

To motivate the QSM condition, we generalize the linear bandit model of Abbasi-Yadkori, Pál, and Szepesvári (2011). Specifically, at each round $t$, the environment generates a context $x(t) \in \mathcal{X} \subseteq \mathbb{R}^d$ (from a possibly adaptive adversary), where $\|x(t)\|_2 \leq S_x$. The learner must then choose an action $a(t) \in \mathcal{A} = \times_{i=1}^n \mathcal{A}_i$, where $\mathcal{A}_i = [K]$. Given a partition $\mathcal{P}$ of $\mathcal{V}$, the learner then receives $|\mathcal{P}|$ noisy observations $\left\{ r_J(t) = \sum_{i \in J} \left\langle x(t), \theta_{i,a_i(t)}^* \right\rangle + \eta_J(t) \right\}_{J \in \mathcal{P}}$, where $\{\theta_{i,j}^* \in \mathbb{R}^d : i \in [n], j \in [K]\}$ are unknown vectors, $\|\theta_{i,j}^*\|_2 \leq S_\theta$, and $\{\eta_J(t)\}_{J \in \mathcal{P}}$ are independent random variables (for every $t$). We assume $\eta_J(t)$ is conditionally $R_J$-subgaussian random noise, such that

$$\mathbb{E}\left[ e^{\lambda \eta_J(t)} \,\Big|\, a_J(1), \ldots, a_J(t), \eta_J(1), \eta_J(t-1) \right] \leq e^{\lambda^2 R_J^2 / 2}.$$

We define the regret at time $T$ by

$$\text{Regret}(T) = \sum_{t=0}^T \sum_{i=1}^n \left[ \left\langle x(t), \theta_{i,a_i^*(t)}^* \right\rangle - \left\langle x(t), \theta_{i,a_i(t)}^* \right\rangle \right],$$

where $a^*(t) \in \arg\max_{a \in \mathcal{A}} \sum_{i=1}^n \left\langle x(t), \theta_{i,a_i}^* \right\rangle$.

Algorithm 1 uses the structured partition under which the QSM condition holds. At every iteration of the algorithm, a least square problem is solved for every $J \in \mathcal{P}$, after which an action is chosen according to an upper confidence defined by

$$\sqrt{\beta_J(t, \delta)} = \lambda^{1/2} |J| S_\theta + R_{\max} \sqrt{d \log\left( \frac{|\mathcal{P}| K^{|J|} (1 + tS_x)/\lambda}{\delta} \right)}.$$

Denote $K_\mathcal{P} = \sum_{J \in \mathcal{P}} K^{|J|}$ and $R_{\max} = \max_{J \in \mathcal{P}} R_J$. Then, we have the following result.

**Theorem 1.** *Assume $\mathbb{E}[r_J] \in [-1, 1]$ for all $J \in \mathcal{P}$. For all $T \geq 0$, with probability at least $1 - \delta$, the regret of Algorithm 1 is bounded by*

$$\text{Regret}(T) \leq 2\sqrt{T} \sqrt{d \log\left( \lambda + \frac{TS_x^2}{Kd} \right) K_\mathcal{P}} \times$$
$$\left( \lambda^{1/2} n S_\theta + R_{\max} \sqrt{d \log\left( \frac{|\mathcal{P}| K^n (1 + tS_x)/\lambda}{\delta} \right)} \right).$$

*This leads to, $\text{Regret}(T) \leq \widetilde{\mathcal{O}}\big( dR_{\max} \sqrt{TK_\mathcal{P}} \big)$.*

The above result achieves regret that is dependent on the maximum subgaussian constant $R_{\max}$ and $\sqrt{K_\mathcal{P}}$. This upper bound is significantly lower than the regret of a naive application of the OFUL algorithm in Abbasi-Yadkori, Pál, and Szepesvári (2011). As the latter doesn't assume the QSM condition, it achieves an exponentially larger regret, $\text{Regret}(T) \leq \widetilde{\mathcal{O}}\big( dR_{\text{tot}} \sqrt{TK^n} \big)$, where $R_{\text{tot}} = \sum_{J \in \mathcal{P}} R_J$. Indeed, whenever $\max_{J \in \mathcal{P}} |J| \ll n$ Algorithm 1 achieves regret which is exponentially smaller in $K$. Particularly, when $\mathcal{P} = \{\{1\}, \{2\}, \ldots \{n\}\}$, we get that $\text{Regret}(T) \leq \widetilde{\mathcal{O}}\big( dR_{\max} \sqrt{TnK} \big)$.

### 3.3 Monotonic Decomposition of Utilities

In Section 3.1 we defined the QSM condition and showed it can significantly improve performance in a linear bandit setting, suggesting its benefits for cooperative MARL. Still, a question arises, when does the QSM condition hold? In this section, we show a sufficient monotonicity assumption under which the QSM condition holds. We formalize this assumption below.

**Assumption 2** (Monotonic Decomposition). *We assume there exists a partition $\mathcal{P}$ of $\mathcal{V}$, utility functions $\{U_i^\pi : S_i \times \mathcal{A}_i \mapsto \mathbb{R}\}_{i \in \mathcal{V}}$, and partition functions $\{F_J^\pi : \mathbb{R}^n \mapsto \mathbb{R}\}_{J \in \mathcal{P}}$ such that for all $J \in \mathcal{P}$,*

$$F_J^\pi(\mathbf{U}(s, a)) = \sum_{i \in J} Q_i^\pi(s, a), \text{ and}$$
$$\nabla_{\mathbf{U}} F_J^\pi \geq \mathbf{0},$$

*where $\mathbf{U}(s, a) := (U_1^\pi(s_1, a_1), \ldots, U_n^\pi(s_n, a_n))^T$.*

**Remark 1.** *The monotonic decomposition assumption generalizes to trajectory-dependent utilities $U_i^\pi : \mathcal{T} \mapsto \mathbb{R}$, such that $F_J^\pi(\mathbf{U}(\tau)) = \sum_{i \in J} Q_i^\pi(s, a)$, where $s, a$ are the final state and action in the trajectory $\tau$.*

A basic setting for which Assumption 2 holds is decoupled MAMDPs. Indeed, for any $\mathcal{M} = (\mathcal{G}, \mathcal{S}, \mathcal{A}, P, r, \gamma)$ such that $\mathcal{E} = \emptyset$ and $\mathcal{M}$ is additively decomposable (see Assumption 1), we have that Assumption 2 holds for any partition $\mathcal{P}$. We refer the reader to the appendix for a proof as well as examples of Assumption 2.

### 3.4 Monotonic Utilities are Sufficient for QSM

Next, we show that monotonic utilities (Assumption 2) are sufficient for the QSM condition. Additionally, we show that, under Assumption 2, local utilities are enough for global $Q$-maximization. This result is closely related to maximization results in previous work (Son et al. 2019; Rashid et al. 2018). See Appendix for proof.

**Theorem 2.** *Suppose Assumption 2 holds for some partition* $\mathcal{P}$. *Then the QSM condition (Definition 1) is satisfied with* $\mathcal{P}$. *Moreover, for any state* $s = (s_1, \ldots, s_n) \in \mathcal{S}$, $\arg\max_{a \in \mathcal{A}} Q^\pi(s, a) = \bigtimes_{i=1}^n \arg\max_{a_i \in \mathcal{A}_i} U_i^\pi(s_i, a_i)$.

Assumption 2 is a generalization of the monotonicity assumption of $Q$-mix (Rashid et al. 2018), which holds when $\mathcal{P} = \{\mathcal{V}\}$. In contrast, when $\mathcal{P} = \{\{1\}, \{2\}, \ldots, \{n\}\}$, each $Q_i$ can be expressed as a function of $\mathbf{U}(\mathbf{s}, \mathbf{a})$, and is monotonic w.r.t to its inputs. The following proposition shows that, if Assumption 2 holds for $\mathcal{P}'$, a refinement of $\mathcal{P}$, then it holds for $\mathcal{P}$ as well. Particularly, this means that if Assumption 2 holds for any $\mathcal{P}$, then the assumption holds for $\mathcal{P} = \{\mathcal{V}\}$.

**Proposition 1.** *Let* $\mathcal{M} = (\mathcal{G}, \mathcal{S}, \mathcal{A}, P, r, \gamma)$, *and let* $\mathcal{P}, \mathcal{P}'$ *be partitions such that* $\mathcal{P}'$ *is a refinement of* $\mathcal{P}$. *If Assumption* 2 *holds for* $\mathcal{P}'$, *then it also holds for* $\mathcal{P}$.

The above proposition suggests a certain trade-off between the refinement of $\mathcal{P}$ and the number of MAMDP's that satisfy Assumption 2. Assumption 2 can thus be viewed as a trade-off between expressibility and speed, as controlled by the refinement of $\mathcal{P}$.

In the next section, we build upon Assumption 2 to construct a scalable value-based MARL algorithm that efficiently leverages local value-partitions and local rewards.

## 4 Local Multi-Agent $Q$-Learning

In this section, we describe a value-based approach that leverages the QSM condition using an application of Assumption 2. Algorithm 2 provides pseudo-code of our method, which we call LOcal Multi-Agent $Q$-learning (LOMA$Q$). We assume local agent rewards are observable during learning (this assumption will be lifted in Section 4.2). Instead of approximating the global $Q$-function, LOMA$Q$ builds upon Assumption 2 to approximate the partition functions $\{F_J\}_{J \in \mathcal{P}}$.

Algorithm 2 receives as input a partition $\mathcal{P}$ and enforces the monotonicity assumption of *Assumption* 2. At every iteration of the algorithm, a greedy action is taken w.r.t. each utility. After an action has been selected, $\{F_J\}_{J \in \mathcal{P}}$ are updated using a bellman update for every $J \in \mathcal{P}$. Finally, in line 10, monotonicity is enforced to ensure Assumption 2 holds. After training is complete, we use the learned utilities

---

**Algorithm 2: LOMA$Q$ with local rewards**

1: **Input:** Partition $\mathcal{P}$ of $\mathcal{V}$, exploration parameter $\epsilon$
2: **Init:** $F_J(\{\mathbf{U}(s', a')\}) = 0$, for all $J \in \mathcal{P}$
3: **for** $t = 1, 2 \ldots$ **do**
4:      Take action $a$
5:      Observe $s'$ and local rewards $\{r_J\}_{J \in \mathcal{P}}$
6:      $a'_{\text{greedy}} \in \left( \arg\max_{a'_i} U_i(s'_i, a'_i) \right)_{i \in \mathcal{V}}$
7:      $a' \leftarrow \begin{cases} \text{random action} & , \text{w.p. } \epsilon \\ a'_{\text{greedy}} & , \text{w.p. } 1 - \epsilon \end{cases}$
8:      **for** $J \in \mathcal{P}$ **do**
9:          $F_J(\mathbf{U}(s, a)) \overset{\alpha_t}{\leftarrow} r_J(s, a) + \gamma F_J(\mathbf{U}(s', a'))$
10:        Project $F_J$ to the set $\{f : \mathbb{R}^n \mapsto \mathbb{R} \text{ s.t. } \nabla f \geq 0\}$
11:      **end for**
12: **end for**

---

$U_i$ for decentralized execution. We note that, due to Theorem 2, choosing the greedy action in line 6 w.r.t. the local utilities is equivalent to acting greedily w.r.t. the global $Q$-function. We refer the reader to the appendix for a discussion regarding the convergence of Algorithm 2.

### 4.1 Practical Implementation of LOMA$Q$

We implement LOMA$Q$ in a deep $Q$-learning framework (Rashid et al. 2018). Specifically, we approximate $F_J^\pi$ for every $J \in \mathcal{P}$, and $U_i^\pi$ for every $i \in \mathcal{V}$ using neural networks with parameters $\theta$. We denote these approximations by $F_J^\theta$ and $U_i^\theta$, respectively. The outputs of $U_i^\theta$ are forwarded as inputs into $F_J^\theta$, i.e., $F_J^\theta(\{U_i^\theta\}_{i=1}^n)$.

Given a mini-batch of tuples $(s, a, r, s')$ sampled from a replay memory, we train the neural networks end-to-end by minimizing the loss

$$L_F(\theta) = \mathbb{E}_{s,a,s'} \left[ \sum_{J \in \mathcal{P}} \left( y_J - F_J^\theta(\{U_i^\theta(s_i, a_i)\}_{i=1}^n) \right)^2 \right], \tag{1}$$

where, $y_J = \sum_{j \in J} r_j + \gamma \max_{a'} \left\{ F_J^\theta(\{U_i^\theta(s'_i, a'_i)\}_{i=1}^n) \right\}$.

Figure 2 depicts the feed-forward architecture for LOMA$Q$. The local agents states $s_i$ are fed into $U_i^\theta$, which outputs a vector of size $\mathcal{A}_i$, representing the utility of every state-action pair $(s_i, a_i)$. The utilities of the chosen actions $a'_i$ are then forwarded as inputs into $F_J^\theta(\{U_i^\theta(s_i, a'_i)\}_{i=1}^n)$. Finally, the outputs of $F_J^\theta$ are trained according to $L_F(\theta)$ in Equation (1).

In practice, every agent $i \in \mathcal{V}$ views a trajectory of local states, represented by $\tau_i$. We use recurrent networks for estimating $U_i^\theta$, and fully-connected networks for $F_J^\theta$. We utilize the graph structure for approximating $F_J$, by redirecting $U_i$ into $F_J$ only if there exists a $j \in J$ such that $i \in N(j)$. We refer the reader to the appendix for an exhaustive overview of specific implementation details.

**Monotonic Regularization** To enforce the monotonicity criterion of Assumption 2, we implement line 10 of Algorithm 2 by regularizing the loss in Equation (1). We propose
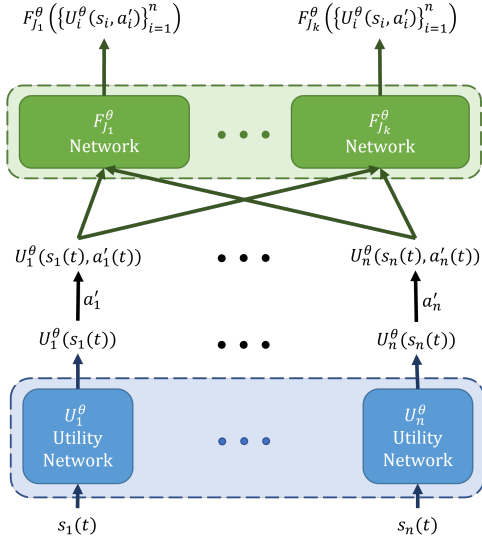
Figure 2: The architecture for the LOMA$Q$ network. The agent's states $s_i$ are fed into the utility networks $U_i^\theta$, which are then forwarded as inputs into $F_J^\theta$.

two such regularizations; namely, using hard and soft projection regularizers.

For hard regularization, we project all parameters $\theta$ to be positive through a Relu activation function, i.e., $\theta \leftarrow \text{Relu}(\theta)$ for all $\theta$ corresponding to $F_J^\theta$. Alternatively, to allow for softer regularization, we penalize Equation (1) by the negative derivatives of $F_J^\theta$ w.r.t. $U_i^\theta$ for every $J \in \mathcal{P}$. That is, $L(\theta) = L_F(\theta) + \lambda \mathcal{R}_{reg}(\theta)$, where $\lambda > 0$, and

$$\mathcal{R}_{reg}(\theta) = \sum_{J \in \mathcal{P}} \text{Relu}(-\nabla_{\mathbf{U}} F_J^\theta).$$

Here, the regularization parameter $\lambda$ reflects a trade-off between efficiency (due to QSM) and accuracy (whenever Assumption 2 does not hold exactly).

## 4.2 Global Reward

While LOMA$Q$ relies on observable local rewards for estimating $F_J^\pi$, they may not always be provided. In this section, we propose a new method for decomposing the global reward function into local reward functions, whenever these are not available.

We assume the global reward signal can be approximately additively decomposed (see Assumption 1). We approximate each local reward $r_i(s_i, a_i)$ using a deep neural network with parameters $\phi$. Our prediction for the global reward is then given by $r_{\text{pred}}^\phi(s, a) = \sum_{i=1}^n r_i^\phi(s_i, a_i)$, which is trained to match the global reward signal $r_{\text{global}}$, by minimizing the loss

$$L_r(\phi) = \mathbb{E}_{s,a}\left[\left(r_{\text{pred}}^\phi(s, a) - r_{\text{global}}(s, a)\right)^2\right]. \quad (2)$$

Training $r_i^\phi(s_i, a_i)$ is done in parallel to LOMA$Q$, where $(s, a, r_{\text{global}})$ are sampled from a replay memory. We refer

the reader to the appendix for an exhaustive overview and further implementation details.

## 4.3 Beyond Additive Decomposition

In certain settings, Assumption 1 may be too restrictive, e.g. when interactions between agents are exhibited in the global reward signal. To overcome this, we consider an alternative decomposition of the reward, where every learned reward function can be dependent on a *group* of agents.

Formally, for any $i \in \mathcal{V}$ we denote by $\mathcal{I}(i)$ the power set of agents in $\{i\} \cup N(i)$. That is,

$$\mathcal{I}(i) = \{I : I \text{ is in the power set of } \{i\} \cup N(i)\}.$$

Next, for every set $I \in \mathcal{I}(i)$ we define a reward function relating to the agents in $I$, $r_I : \mathcal{S}_I \times \mathcal{A}_I \mapsto \mathbb{R}$. Finally, we define the reward of agent $i \in \mathcal{V}$ by

$$r_i(s, a) = \sum_{I \in \mathcal{I}(i)} \frac{1}{|I|} r_I(s_I, a_I), \quad (3)$$

where here, every reward $r_I$ in the summand is normalized according to the cardinality of $I$. Notice that this decomposition is a generalization of Assumption 1. Indeed, Equation (3) coincides with Assumption 1 whenever $\mathcal{E} = \emptyset$.

The reward decomposition in Equation (3) creates a hierarchy for every agent $i$, as every local reward $r_i$ is comprised of multiple learned reward functions $\{r_I(s_I, a_I)\}_{I \in \mathcal{I}(i)}$ which have less effect on agent $i$ as $|I|$ increases.

In most cases, the number of local reward functions is exponential in $N(i)$, rendering large decompositions infeasible. Moreover, as $|I|$ increases, the learned rewards become dependent on more agents, reducing their effectiveness (due to normalization in $|I|$). We therefore focus on learning reward functions of small cardinality in $|I|$. We enforce this in practice using a regularization term that is dependent on $|I|$. Specifically, we regularize the loss in Equation (2) by

$$\mathcal{R}_{\text{reg}}(\phi) = \sum_{I \in \mathcal{I}} w(|I|) \times |r_I^\phi(s_I, a_I)|,$$

where $w(|I|)$ are weights that grow proportionally to $|I|$, penalizing $r_I^\phi(s_I, a_I)$ as $|I|$ increases. This regularization reflects a trade-off between the overall accuracy of the learned rewards and the complexity of the reward decomposition.

# 5 Experiments

In this section we test the performance of LOMA$Q$ and compare it to previous MARL approaches on two large-scale multi agent tasks.

## 5.1 Environments

We tested our algorithm on two environments, Coupled-Multi-Cart-Pole and Bounded-Cooperative-Navigation. Both environments include minor modifications of the well-known Cart-Pole (Brockman et al. 2016) and Cooperative-Navigation (Lowe et al. 2017) environments.

The Coupled-Multi-Cart-Pole consists of $n$ cartpoles, residing on the 1d axis. Each cart is viewed as an agent, controlled by applying a force of $\pm 1$. Every pair of neighboring

Figure 3: The Coupled-Multi-Cart-Pole environment with 3 cartpoles. The right-most cartpole has fallen (marked in red). The global reward for this timestep is +2.



(a) Particles in purple, landmarks in red, regions in gray    (b) Interaction graph based on region overlap

Figure 4: The Bounded-Cooperative-Navigation environment with 16 agents and circular regions.

carts is connected by a spring. Every cart receives a local reward of $+1$ for every timestep that the pole is upright. The global reward for the environment is the total number of cartpoles that are currently upright. The dependency graph for this environment can be modeled as a line graph, where every cartpole has two neighbors excluding the edges which only have one. Figure 3 depicts this environment for three cartpoles.

The Bounded-Cooperative-Navigation consists of $n$ agents (particles) and $n$ landmarks. Agents must strive to cooperatively cover as many landmarks as possible. In this environment, particles aren't able to move freely in 2d space. Every particle is bound to a fixed distance from its starting position and is thereby restricted to a certain region. This restriction resembles a simplified food-delivery service, where the landmarks represent customers and the agents represent delivery people. Consequently, not all particles interact with each other directly, since direct interactions only occur when two particles are in the same location. This induces a dependency graph for which every two particles are neighbors in the graph if and only if their regions overlap. The environment rewards $+1$ for every landmark that is covered by a particle at a certain timestep. Figure 4 shows a conceptual visualization of the task.

### 5.2 Comparative Experiments

**Scalability**   We tested LOMA$Q$ on both cooperative environments with $n = 15$ agents in two setups; namely, with and without access to a local reward signal. We denote these by LOMA$Q$ and LOMA$Q$+RD, respectively. In both cases, we used the refined partition $\mathcal{P} = \{\{1\}, \{2\}, \dots, \{n\}\}$.

We compared LOMA$Q$ to a wide range of contemporary cooperative methods. In addition, we compared LOMA$Q$ to two versions of IQL (Tan 1993), trained with environment local rewards and global rewards, which we denote by IQL-local and IQL, respectively.

Figure 5 depicts the results of the Coupled-Multi-Cart-Pole and Bounded-Cooperative-Navigation environments. It is evident that both versions of LOMA$Q$ significantly outperform all of the compared methods in both performance and convergence speed. We note that LOMA$Q$+RD converges to LOMA$Q$'s policy, with a slight delay due to the time taken to learn the reward decomposition.

In both environments various cooperative methods exhibit slow learning compared to LOMA$Q$, due to the use of global rewards. Additionally, while IQL-local does learn quickly
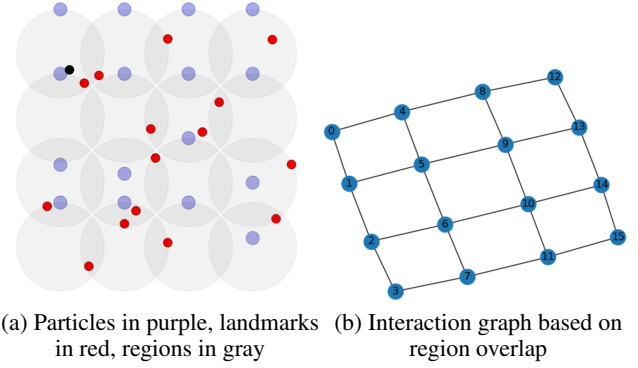
(primarily due to the use of local rewards), it converged to a sub-optimal solution. This occurs as IQL acts greedily w.r.t its local rewards. In contrast, LOMA$Q$ incentivizes cooperation, enabling both fast convergence as well as improved performance.

**Reward Decomposition**   We visualize multiple reward decompositions for Bounded-Cooperative-Navigation. We run our decomposition method with a global reward signal, for $n = 2$ agents and a single landmark. If both agents are on the landmark at the same time, the global reward remains 1. We plot the learned reward functions as a function of Agent 1 and Agent 2's distance from the landmark, which we denote by $\Delta x$. These results are depicted in Figure 6.

The first row in Figure 6 assumes a decomposition according to Assumption 1. Assumption 1 does not hold for this setup, since the reward function is dependent on both agents when they share a landmark. The approximated reward is overly optimistic and wrongly rewards $+2$ when the landmark is shared. The second row approximates a decomposition that allows $|I| \leq 2$ with no regularization $\lambda = 0$. In this case, the global reward is approximated correctly, and the local reward functions $r^{\phi}_{\{i\}} = 0$. The third row visualizes a decomposition with regularization $w(|I| = 2) = 1, \lambda = 0.0001$. In this case, the local reward functions $r^{\phi}_{\{i\}}$ convey information for each agent $i$, and $r^{\phi}_{\{1,2\}}$ conveys information regarding their joint dynamic.

## 6   Related Work

**Graph Based MARL.**   The underlying structure of the team of agents in the environment can often be modeled using a graph topology. Jiang et al. (2019) propose DGN - a MARL algorithm based on the graph convolutional network (GCN) architecture which assumes centralized execution and homogeneous agents. Naderializadeh et al. (2020) propose GraphMIX for CTDE, that uses global rewards for learning. Qu et al. (2020) propose Scalable Actor-Critic - An Actor-Critic approach for the discrete space case which utilizes a dependency graph with theoretical guarantees.
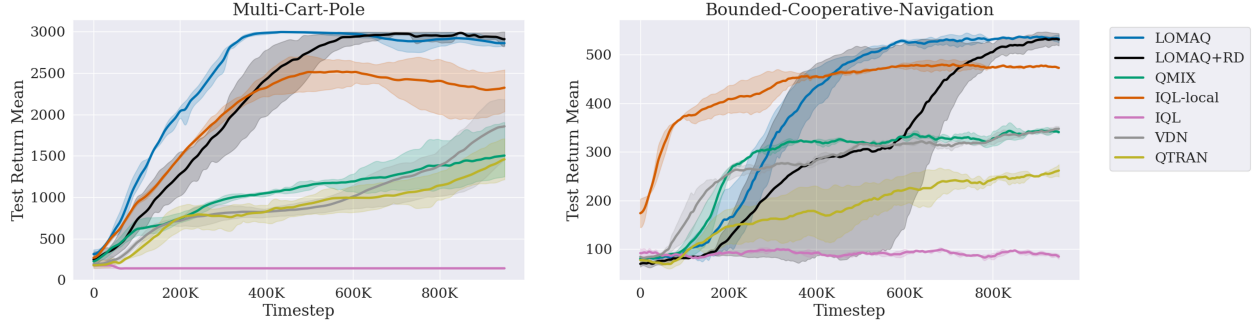
Figure 5: Test returns for the Coupled-Multi-Cart-Pole environment and for the Bounded-Cooperative-Navigation environment.
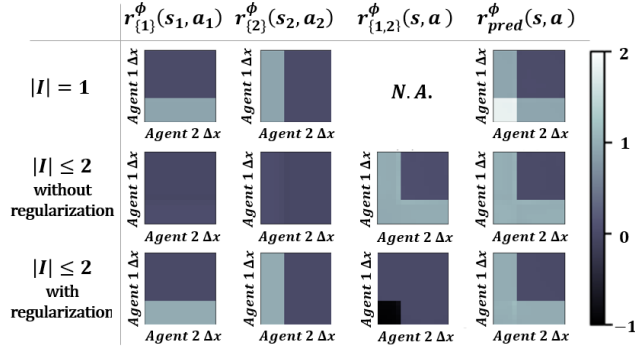


Figure 6: Visualization of learned reward functions $r_I^\phi$ for different decompositions in Bounded-Cooperative-Navigation for $n = 2$ agents and a single landmark. We plot the learned reward functions as a function of Agent 1 and Agent 2's distance from the landmark, denoted by $\Delta x$. The first row assumes $|I| \leq 1$, the second row assumes $|I| \leq 2$ with no regularization, and the third row adds regularization.

**Cooperative MARL.** Our approach enhances the popular value decomposition family (Son et al. 2019; Wang et al. 2020; Rashid et al. 2020), which consider a cooperative multi-agent problem in which each agent observes its own state and action history. Sunehag et al. (2018) propose VDN for decomposing the value function into a sum of utility functions. Rashid et al. (2018) offer $Q$-mix which generalizes this concept, by decomposing the value function into a monotonic function of individual utility functions. All of these approaches implicitly measure the impact of every agent on the observed global reward, whereas we propose to combine this line of work with an explicit approach for credit assignment using local rewards.

**Credit Assignment.** Various approaches have attempted to tackle the credit assignment problem. A common approach for credit assignment is by estimating the individual $Q$ functions $Q_i$ directly, which are often substantially simpler and significantly easier to learn than $Q$ (Qu et al. 2020; Kok and Vlassis 2004; Russell and Zimdars 2003; van Seijen et al. 2017; Juozapaitis et al. 2019). Our work extends this line of work for value-based CTDE, and focuses on re-

ward decompositions that expedite learning alongside global cooperation.

**Reward Decomposition.** Multiple works recognize the benefits of local rewards and attempt to learn them in settings where only a global reward signal is provided. $RD^2$ (Lin et al. 2020) learns a reward decomposition with minimally-dependent features for factored-state MDP setting. Our method can be seen as an extension of $RD^2$ for MARL, where the action is also factored.

**Large Action Spaces.** Finally, our work is related to work on large and combinatorial action spaces. From action elimination (Zahavy et al. 2018), to action embeddings (Tennenholtz and Mannor 2019; Chandak et al. 2019), through action redundancy (Baram, Tennenholtz, and Mannor 2021), our work can be viewed as an additional method for reducing the effective dimensionality of the problem.

# 7 Conclusion and Future Work

In this work we tackled the credit assignment problem of cooperative MARL through local, partition based value functions. We used the QSM condition and a monotonic decomposition of utilities to construct a value-based approach, effectively reducing the problem to simpler ones. We showed that local rewards are highly beneficial, both when provided as well as learned implicitly from a global reward. These greatly improved overall performance and convergence speed, suggesting that local structures can be efficiently used to improve MARL algorithms.

In this work we have assumed that an underlying, static dependency graph $\mathcal{G}$ is provided during training. In many cases, these assumptions are limiting. We look to further generalize our method by learning such dynamic dependencies between agents through interaction with the environment. In addition, our work has assumed that $Assumption$ 2 holds for some partition $\mathcal{P}$ and local reward decomposition $\{r_i\}$. We look to generalize our algorithm to automatically identify effective decompositions.

# References

Abbasi-Yadkori, Y.; Pál, D.; and Szepesvári, C. 2011. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24: 2312–2320.

Baram, N.; Tennenholtz, G.; and Mannor, S. 2021. Action Redundancy in Reinforcement Learning. *arXiv preprint arXiv:2102.11329*.

Brockman, G.; Cheung, V.; Pettersson, L.; Schneider, J.; Schulman, J.; Tang, J.; and Zaremba, W. 2016. OpenAI Gym.

Chandak, Y.; Theocharous, G.; Kostas, J.; Jordan, S.; and Thomas, P. 2019. Learning action representations for reinforcement learning. In *International Conference on Machine Learning*, 941–950. PMLR.

Chu, T.; Wang, J.; Codecà, L.; and Li, Z. 2019. Multi-agent deep reinforcement learning for large-scale traffic signal control. *IEEE Transactions on Intelligent Transportation Systems*, 21(3): 1086–1095.

Jiang, J.; Dun, C.; Huang, T.; and Lu, Z. 2019. Graph Convolutional Reinforcement Learning. In *International Conference on Learning Representations*.

Juozapaitis, Z.; Koul, A.; Fern, A.; Erwig, M.; and Doshi-Velez, F. 2019. Explainable reinforcement learning via reward decomposition. In *IJCAI/ECAI Workshop on Explainable Artificial Intelligence*.

Kok, J. R.; and Vlassis, N. 2004. Sparse cooperative Q-learning. In *Proceedings of the twenty-first international conference on Machine learning*, 61.

Kraemer, L.; and Banerjee, B. 2016. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing*, 190: 82–94.

Lin, Z.; Yang, D.; Zhao, L.; Qin, T.; Yang, G.; and Liu, T.-Y. 2020. RD2: Reward Decomposition with Representation Decomposition. *Advances in Neural Information Processing Systems*, 33.

Lowe, R.; WU, Y.; Tamar, A.; Harb, J.; Pieter Abbeel, O.; and Mordatch, I. 2017. Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments. *Advances in Neural Information Processing Systems*, 30: 6379–6390.

Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; et al. 2015. Human-level control through deep reinforcement learning. *nature*, 518(7540): 529–533.

Naderializadeh, N.; Hung, F. H.; Soleyman, S.; and Khosla, D. 2020. Graph Convolutional Value Decomposition in Multi-Agent Reinforcement Learning. *CoRR*, abs/2010.04740.

Qu, G.; Lin, Y.; Wierman, A.; and Li, N. 2020. Scalable Multi-Agent Reinforcement Learning for Networked Systems with Average Reward. In *Thirty-fourth Conference on Neural Information Processing Systems*.

Rashid, T.; Farquhar, G.; Peng, B.; and Whiteson, S. 2020. Weighted QMIX: Expanding Monotonic Value Function Factorisation for Deep Multi-Agent Reinforcement Learning. *Advances in Neural Information Processing Systems*, 33.

Rashid, T.; Samvelyan, M.; Schroeder, C.; Farquhar, G.; Foerster, J.; and Whiteson, S. 2018. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, 4295–4304. PMLR.

Russell, S. J.; and Zimdars, A. 2003. Q-decomposition for reinforcement learning agents. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, 656–663.

Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.; Van Den Driessche, G.; Schrittwieser, J.; Antonoglou, I.; Panneershelvam, V.; Lanctot, M.; et al. 2016. Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587): 484–489.

Son, K.; Kim, D.; Kang, W. J.; Hostallero, D. E.; and Yi, Y. 2019. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*, 5887–5896. PMLR.

Sunehag, P.; Lever, G.; Gruslys, A.; Czarnecki, W. M.; Zambaldi, V. F.; Jaderberg, M.; Lanctot, M.; Sonnerat, N.; Leibo, J. Z.; Tuyls, K.; et al. 2018. Value-Decomposition Networks For Cooperative Multi-Agent Learning Based On Team Reward. In *AAMAS*.

Tan, M. 1993. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the tenth international conference on machine learning*, 330–337.

Tennenholtz, G.; and Mannor, S. 2019. The natural language of actions. In *International Conference on Machine Learning*, 6196–6205. PMLR.

van Seijen, H.; Fatemi, M.; Romoff, J.; Laroche, R.; Barnes, T.; and Tsang, J. 2017. Hybrid reward architecture for reinforcement learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 5398–5408.

Vinyals, O.; Babuschkin, I.; Czarnecki, W. M.; Mathieu, M.; Dudzik, A.; Chung, J.; Choi, D. H.; Powell, R.; Ewalds, T.; Georgiev, P.; et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782): 350–354.

Wang, J.; Ren, Z.; Liu, T.; Yu, Y.; and Zhang, C. 2020. QPLEX: Duplex Dueling Multi-Agent Q-Learning. In *International Conference on Learning Representations*.

Zahavy, T.; Haroush, M.; Merlis, N.; Mankowitz, D. J.; and Mannor, S. 2018. Learn What Not to Learn: Action Elimination with Deep Reinforcement Learning. *Advances in Neural Information Processing Systems*, 31: 3562–3573.