

ROMA: Multi-Agent Reinforcement Learning with Emergent Roles

Tonghan Wang¹ Heng Dong¹ Victor Lesser² Chongjie Zhang¹

Abstract

The role concept provides a useful tool to design and understand complex multi-agent systems, which allows agents with a similar role to share similar behaviors. However, existing role-based methods use prior domain knowledge and predefine role structures and behaviors. In contrast, multi-agent reinforcement learning (MARL) provides flexibility and adaptability, but less efficiency in complex tasks. In this paper, we synergize these two paradigms and propose a role-oriented MARL framework (ROMA). In this framework, roles are emergent, and agents with similar roles tend to share their learning and to be specialized on certain sub-tasks. To this end, we construct a stochastic role embedding space by introducing two novel regularizers and conditioning individual policies on roles. Experiments show that our method can learn specialized, dynamic, and identifiable roles, which help our method push forward the state of the art on the StarCraft II micromanagement benchmark. Demonstrative videos are available at <https://sites.google.com/view/romar1/>.

1. Introduction

Many real-world systems can be modeled as multi-agent systems (MAS), such as autonomous vehicle teams (Cao et al., 2012), intelligent warehouse systems (Nowé et al., 2012), and sensor networks (Zhang & Lesser, 2011). Cooperative multi-agent reinforcement learning (MARL) provides a promising approach to developing these systems, allowing agents to deal with uncertainty and adapt to the dynamics of an environment. In recent years, cooperative MARL has achieved prominent progress, and many deep methods have been proposed (Foerster et al., 2018; Sunehag et al., 2018;

¹IIIS, Tsinghua University, Beijing, China ²University of Massachusetts, Amherst, USA. Correspondence to: Tonghan Wang <tonghanwang1996@gmail.com>.

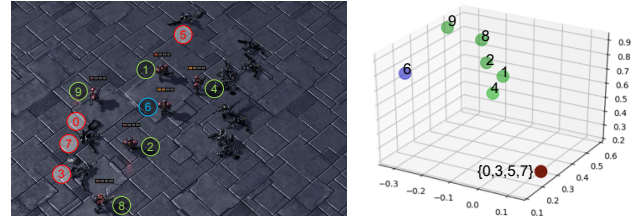


Figure 1. Visualization of our learned role representations at a timestep. The blue agent has the maximum health, while the red ones are dead. The corresponding policy is that agent 6 moves towards enemies to take on more firepower, so that more seriously injured agents are protected. Roles can change adaptively and will aggregate according to responsibilities that are compatible with individual characteristics, such as location, agent type, health, etc.

Rashid et al., 2018; Son et al., 2019; Vinyals et al., 2019; Wang et al., 2020b; Baker et al., 2020).

In order to achieve scalability, these deep MARL methods adopt a simple mechanism that all agents share and learn a decentralized value or policy network. However, such simple sharing is often not effective for many complex multi-agent tasks. For example, in Adam Smith’s Pin Factory, workers must complete up to eighteen different tasks to create one pin (Smith, 1937). In this case, it is a heavy burden for a single shared policy to represent and learn all required skills. On the other hand, it is also unnecessary for each agent to use a distinct policy network, which leads to high learning complexity because some agents often perform similar sub-tasks from time to time. The question is how we can give full play to agents’ specialization and dynamic sharing for improving learning efficiency.

A natural concept that comes to mind is the *role*. A role is a comprehensive pattern of behavior, often specialized in some tasks. Agents with similar roles will show similar behaviors, and thus can share their experiences to improve performance. The role theory has been widely studied in economics, sociology, and organization theory. Researchers have also introduced the concept of role into MAS (Becht et al., 1999; Stone & Veloso, 1999; Depke et al., 2001; Ferber et al., 2003; Odell et al., 2004; Bonjean et al., 2014; Lhaksmana et al., 2018). In these role-based frameworks, the complexity of agent design is reduced via task decomposition by defining roles associated with responsibilities

made up of a set of sub-tasks, so that the policy search space is effectively decomposed (Zhu & Zhou, 2008). However, these works exploit prior domain knowledge to decompose tasks and predefine the responsibilities of each role, which prevents role-based MAS from being dynamic and adaptive to uncertain environments.

To leverage the benefits of both role-based and learning methods, in this paper, we propose a role-oriented multi-agent reinforcement learning framework (ROMA). This framework implicitly introduces the role concept into MARL, which serves as an intermediary to enable agents with similar responsibilities to share their learning. We achieve this by ensuring that agents with similar roles have both similar policies and responsibilities. To establish the connection between roles and decentralized policies, ROMA conditions agents’ policies on individual roles, which are stochastic latent variables determined by agents’ local observations. To associate roles with responsibilities, we introduce two regularizers to enable roles to be identifiable by behaviors and specialized in certain sub-tasks. We show how well-formed role representations can be learned via optimizing tractable variational estimations of the proposed regularizers. In this way, our method synergizes role-based and learning methods while avoiding their individual shortcomings – we provide a flexible and general-purpose mechanism that promotes the *emergence* and *specialization* of roles, which in turn provides an adaptive learning sharing mechanism for efficient multi-agent policy learning.

We test our method on StarCraft II¹ micromanagement environments (Vinyals et al., 2017; Samvelyan et al., 2019). Results show that our method significantly pushes forward the state of the art of MARL algorithms, by virtue of the adaptive policy sharing among agents with similar roles. Visualization of the role representations in both homogeneous and heterogeneous agent teams demonstrates that the learned roles can adapt automatically in dynamic environments, and that agents with similar responsibilities have similar roles. In addition, the emergence and evolution process of roles is shown, highlighting the connection between role-driven sub-task specialization and improvement of team efficiency in our framework. These results provide a new perspective in understanding and promoting the emergence of cooperation among agents.

2. Background

In our work, we consider a fully cooperative multi-agent task that can be modelled by a Dec-POMDP (Oliehoek et al., 2016) $G = \langle I, S, A, P, R, \Omega, O, n, \gamma \rangle$, where A is the finite action set, I is the finite set of n agents, $\gamma \in [0, 1)$ is the discount factor, and $s \in S$ is the true state of the environ-

ment. We consider partially observable settings and agent i only has access to an observation $o_i \in \Omega$ drawn according to the observation function $O(s, i)$. Each agent has a history $\tau_i \in T \equiv (\Omega \times A)^*$. At each timestep, each agent i selects an action $a_i \in A$, forming a joint action $\mathbf{a} \in A^n$, leading to next state s' according to the transition function $P(s'|s, \mathbf{a})$ and a shared reward $r = R(s, \mathbf{a})$ for each agent. The joint policy π induces a joint action-value function: $Q_{tot}^\pi(s, \mathbf{a}) = \mathbb{E}_{s_0, \infty, \mathbf{a}_0, \infty} [\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s, \mathbf{a}_0 = \mathbf{a}, \pi]$.

To effectively learn policies for agents, the paradigm of centralized training with decentralized execution (CTDE) (Foster et al., 2016; 2018; Wang et al., 2020a) has recently attracted attention from deep MARL to deal with non-stationarity while learning decentralized policies. One of the promising ways to exploit the CTDE paradigm is value function decomposition (Sunehag et al., 2018; Rashid et al., 2018; Son et al., 2019; Wang et al., 2020b), which learns a decentralized utility function for each agent and uses a mixing network to combine these local utilities into a global action value. To achieve learning scalability, existing CTDE methods typically learn a shared local value or policy network for agents. However, this simple sharing mechanism is often not sufficient for learning complex tasks, where diverse responsibilities or skills are required to achieve goals. In this paper, we develop a novel role-based MARL framework to address this challenge. This framework achieves efficient shared learning while allowing agents to learn sufficiently diverse skills.

3. Method

In this section, we will present a novel role-oriented MARL framework (ROMA) that introduces the role concept into MARL and enables adaptive shared learning among agents. ROMA adopts the CTDE paradigm. As shown in Fig. 2, it learns local Q-value functions for agents, which are fed into a mixing network to compute a global TD loss for centralized training. During the execution, the mixing network will be removed, and each agent will act based on its local policy derived from its value function. Agents’ value functions or policies are dependent on their roles, each of which is responsible for performing similar automatically identified sub-tasks. To enable efficient and effective shared learning among agents with similar behaviors, ROMA will automatically learn roles that are:

- i) **Dynamic:** An agent’s role can automatically adapt to the dynamics of the environment;
- ii) **Identifiable:** The role of an agent contains enough information about its behaviors;
- iii) **Specialized:** Agents with similar roles are expected to specialize in similar sub-tasks.

¹StarCraft II are trademarks of Blizzard Entertainment™.

Formally, each agent i has a local utility function (or an individual policy), whose parameters θ_i are conditioned on its role ρ_i . To learn roles with desired properties, we encode roles in a stochastic embedding space, and the role of agent i , ρ_i , is drawn from a multivariate Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}_{\rho_i}, \boldsymbol{\sigma}_{\rho_i})$. To enable the **dynamic** property, ROMA conditions an agent’s role on its local observations, and uses a trainable neural network f to learn the parameters of the Gaussian distribution of the role:

$$\begin{aligned} (\boldsymbol{\mu}_{\rho_i}, \boldsymbol{\sigma}_{\rho_i}) &= f(o_i; \theta_\rho), \\ \rho_i &\sim \mathcal{N}(\boldsymbol{\mu}_{\rho_i}, \boldsymbol{\sigma}_{\rho_i}), \end{aligned} \quad (1)$$

where θ_ρ are parameters of f . The sampled role ρ_i is then fed into a hyper-network $g(\rho_i; \theta_h)$ parameterized by θ_h to generate the parameters for the individual policy, θ_i . We call f the *role encoder* and g the *role decoder*. In the next two sub-sections, we will describe two regularizers for learning identifiable and specialized roles.

3.1. Identifiable Roles

Introducing latent role embedding and conditioning individual policies on this embedding does not automatically generate roles with desired properties. Intuitively, conditioning roles on local observations enables roles to be responsive to the changes in the environment. This design enables ROMA to be adaptive to dynamic environments but may cause roles to change quickly, making learning unstable. For addressing this problem, we expect roles to be temporally stable. To this end, we propose to learn roles that are identifiable by agents’ long term behaviors, which can be achieved by maximizing $I(\tau_i; \rho_i | o_i)$, the conditional mutual information between the individual trajectory and the role given the current observation.

However, estimating and maximizing mutual information is often intractable. Drawing inspiration from the literature of variational inference (Wainwright et al., 2008; Alemi et al., 2017), we introduce a variational posterior estimator to derive a tractable lower bound for the mutual information objective (the proof is deferred to Appendix A.1):

$$I(\rho_i^t; \tau_i^{t-1} | o_i^t) \geq \mathbb{E}_{\rho_i^t, \tau_i^{t-1}, o_i^t} \left[\log \frac{q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)}{p(\rho_i^t | o_i^t)} \right], \quad (2)$$

where $\tau_i^{t-1} = (o_i^0, a_i^0, \dots, o_i^{t-1}, a_i^{t-1})$, q_ξ is the variational estimator parameterised with ξ . For q_ξ , we use a GRU (Cho et al., 2014) to encode an agent’s history of observations and actions, and call it the *trajectory encoder*. The lower bound in Eq. 2 can be further rewritten as a loss function to be minimized:

$$\mathcal{L}_I(\theta_\rho, \xi) = \mathbb{E}_{(\tau_i^{t-1}, o_i^t) \sim \mathcal{D}} \left[D_{\text{KL}}[p(\rho_i^t | o_i^t) \| q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)] \right], \quad (3)$$

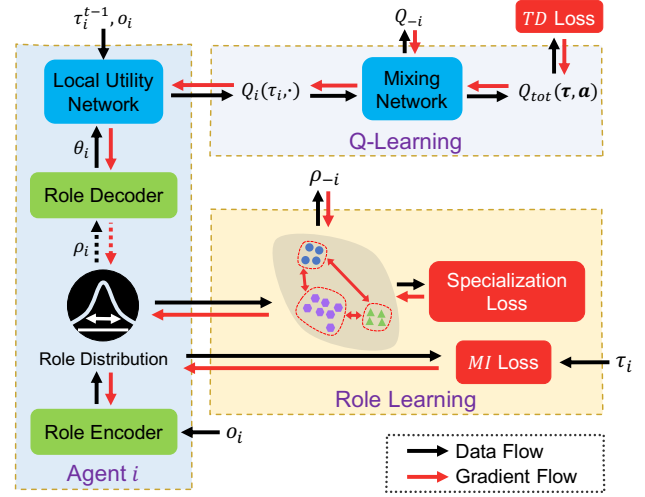


Figure 2. Schematics of our approach. The role encoder generates a role embedding distribution, from which a role is sampled and serves as the input to the role decoder. The role decoder generates the parameters of the local utility network. Local utilities are fed into a mixing network to get an estimation of the global action value. We propose two learning objectives to learn specialized and identifiable roles. The framework can be trained in an end-to-end manner.

where \mathcal{D} is a replay buffer, and $D_{\text{KL}}[\cdot \| \cdot]$ is the KL divergence operator. The detailed derivation can be found in Appendix A.1.

3.2. Specialized Roles

The formulation so far does not promote sub-task specialization, which is the critical component to share learning and improve efficiency in multi-agent systems. Minimizing \mathcal{L}_I enables roles to contain enough information about long-term behaviors but does not explicitly ensure agents with similar behaviors to have similar role embeddings.

For learning specialized roles, we define another role-learning regularizer. Intuitively, to encourage sub-task specialization, for any two agents, we expect that either they have similar roles or they have quite different behaviors. However, it is usually unclear which agents will have similar roles during the process of role emergence, and the similarity between behaviors is not straightforward to measure.

Since roles have enough information about the behaviors (achieved by minimizing \mathcal{L}_I), to encourage two agents i and j to have similar roles, we can maximize $I(\rho_i; \tau_j)$, the mutual information between the role of agent i and the trajectory of agent j . However, we do not know which agents will have similar roles, and directly optimizing this objective for all pairs of agents will result in all agents having the same role, and, correspondingly, the same policy, which

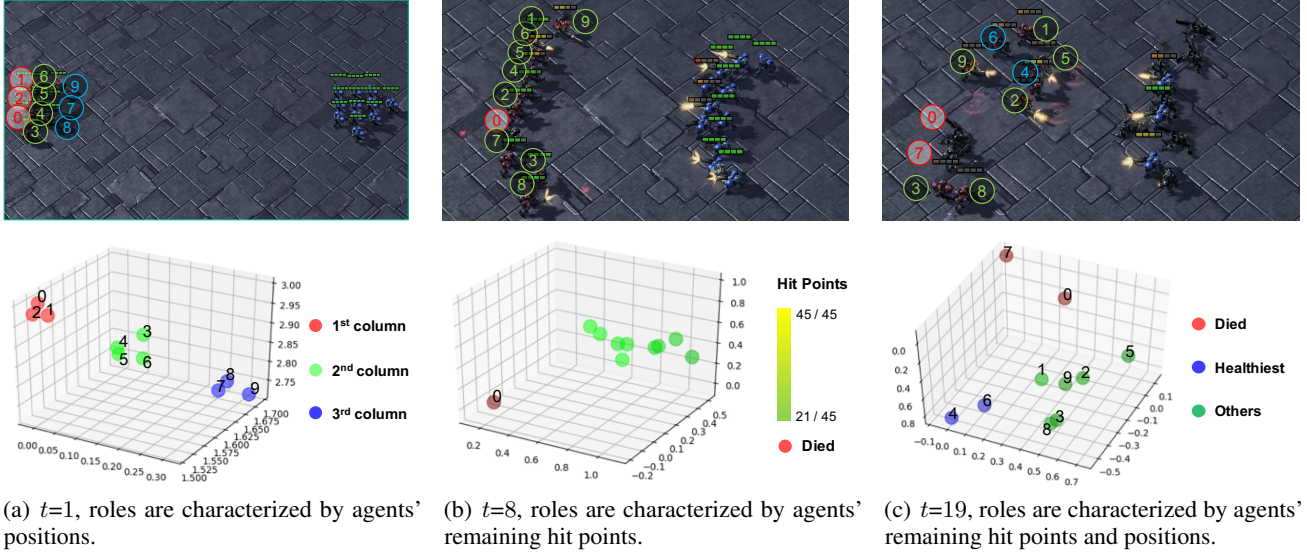


Figure 3. Dynamic role adaptation during an episode (means of the role distributions, μ_{ρ_i} , are shown, without using any dimensionality reduction techniques). The *role encoder* learns to focus on different parts of observations according to the automatically discovered demands of the task. The role-induced strategy helps (a) quickly form the offensive arc when $t=1$; (b) protect injured agents when $t=8$; (c) protect dying agents and alternate fire when $t=19$.

will limit system performance. To settle this issue, we introduce a dissimilarity model $d_\phi : \mathcal{T} \times \mathcal{T} \rightarrow \mathbb{R}$, a trainable neural network taking two trajectories as input, and seek to maximize $I(\rho_i; \tau_j) + d_\phi(\tau_i, \tau_j)$ while minimizing the number of non-zero elements in the matrix $D_\phi = (d_{ij})$. Here, $d_{ij} = d_\phi(\tau_i, \tau_j)$ is the estimated dissimilarity between trajectories of agent i and j . Such formulation makes sure that dissimilarity d is high only when mutual information I is low, so that the set of learned roles is compact but diverse, which help solve the given task efficiently. Formally, the following learning objective encourages sub-task specialization:

$$\begin{aligned} & \underset{\theta_\rho, \xi, \phi}{\text{minimize}} \quad \|D_\phi^t\|_{2,0} & (4) \\ & \text{subject to} \quad I(\rho_i^t; \tau_j^{t-1} | o_j^t) + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}) > U, \forall i \neq j, \end{aligned}$$

where U controls the compactness of the role representation. In practice, we separately carry out min-max normalization on I and d to scale their values to $[0, 1]$ and set U to 1. Relaxing the matrix norm $\|\cdot\|_{2,0}$ with the Frobenius norm, we can get the optimization objective for minimizing:

$$\|D_\phi^t\|_F - \sum_{i \neq j} \min\{I(\rho_i^t; \tau_j^{t-1} | o_j^t) + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\}, \quad (5)$$

However, as estimating and optimizing the mutual information term are intractable, we use the variational posterior estimator introduced in Sec. 3.1 to construct an upper bound, serving as the second regularizer of ROMA:

$$\mathcal{L}_D(\theta_\rho, \phi, \xi) = \mathbb{E}_{(\tau^{t-1}, o^t) \sim \mathcal{D}, \rho^t \sim p(\rho^t | o^t)} [\|D_\phi^t\|_F] \quad (6)$$

$$- \sum_{i \neq j} \min\{q_\xi(\rho_i^t | \tau_j^{t-1}, o_j^t) + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\}$$

where \mathcal{D} is the replay buffer, τ^{t-1} is the joint trajectory, o^t is the joint observation, and $\rho^t = \langle \rho_1^t, \rho_2^t, \dots, \rho_n^t \rangle$. A detailed derivation can be found in Appendix A.2.

3.3. Overall Optimization Objective

We have introduced optimization objectives for learning roles to be identifiable and specialized. Apart from these regularizers, all the parameters in the framework are updated by gradients induced by the standard TD loss of reinforcement learning. As shown in Fig. 2, to compute the global TD loss, individual utilities are fed into a mixing network whose output is the estimation of global action-value Q_{tot} . In this paper, our ROMA implementation uses the mixing network introduced by QMIX (Rashid et al., 2018) (see Appendix D) for its monotonic approximation, but it can be easily replaced by other mixing methods. The parameters of the mixing network are conditioned on the global state s and are generated by a hyper-net parameterized by θ_m . Therefore, the final learning objective of ROMA is:

$$\mathcal{L}(\theta) = \mathcal{L}_{TD}(\theta) + \lambda_I \mathcal{L}_I(\theta_\rho, \xi) + \lambda_D \mathcal{L}_D(\theta_\rho, \xi, \phi), \quad (7)$$

where $\theta = (\theta_\rho, \xi, \phi, \theta_h, \theta_m)$, λ_I and λ_D are scaling factors, and $\mathcal{L}_{TD}(\theta) = [r + \gamma \max_{\mathbf{a}'} Q_{tot}(s', \mathbf{a}'; \theta^-) - Q_{tot}(s, \mathbf{a}; \theta)]^2$ (θ^- are the parameters of a periodically updated target network). In our centralized training with decentralized execution framework, only the role encoder, the role decoder, and the individual utility networks are used when execution.

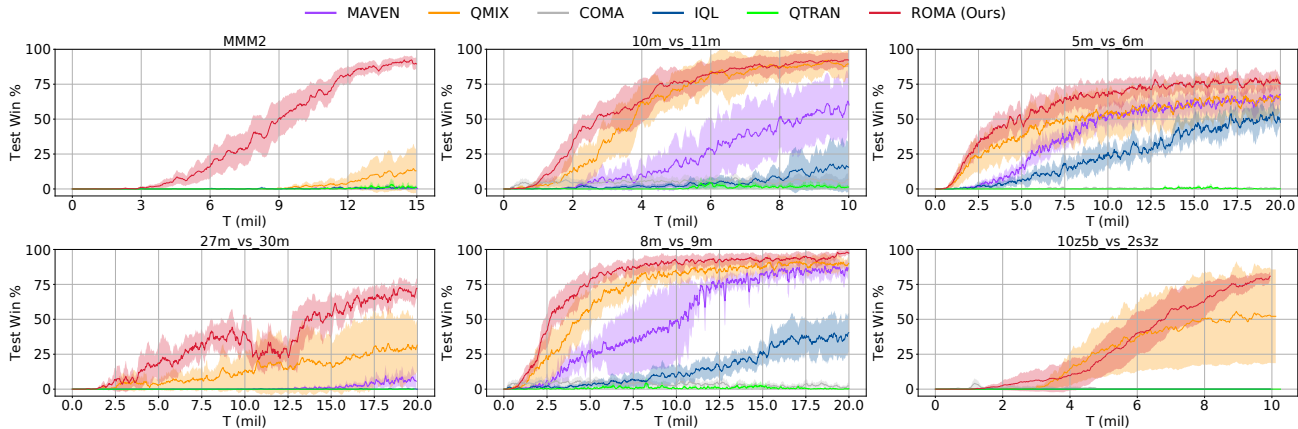


Figure 4. Comparison of our method against baseline algorithms. Results for more maps can be found in Appendix C.1.

4. Related Works

The emergence of role has been documented in many natural systems, such as bees (Jeanson et al., 2005), ants (Gordon, 1996), and humans (Butler, 2012). In these systems, the role is closely related to the division of labor and is crucial to the improvement of labor efficiency. Many multi-agent systems are inspired by these natural systems. They decompose the task, make agents with the same role specialize in certain sub-tasks, and thus reduce the design complexity (Wooldridge et al., 2000; Omicini, 2000; Padgham & Winikoff, 2002; Pavón & Gómez-Sanz, 2003; Cossentino et al., 2005; Zhu & Zhou, 2008; Spanoudakis & Moraitis, 2010; DeLoach & Garcia-Ojeda, 2010; Bonjean et al., 2014). These methodologies are designed for tasks with a clear structure, such as software engineering (Bresciani et al., 2004). Therefore, they tend to use predefined roles and associated responsibilities (Lhaksmana et al., 2018). In contrast, we focus on how to implicitly introduce the concept of roles into general multi-agent sequential decision making under dynamic and uncertain environments.

Deep multi-agent reinforcement learning has witnessed vigorous progress in recent years. COMA (Foerster et al., 2018), MADDPG (Lowe et al., 2017), PR2 (Wen et al., 2019), and MAAC (Iqbal & Sha, 2019) explore multi-agent policy gradients. Another line of research focuses on value-based multi-agent RL, and value-function factorization is the most popular method. VDN (Sunehag et al., 2018), QMIX (Rashid et al., 2018), and QTRAN (Son et al., 2019) have progressively enlarged the family of functions that can be represented by the mixing network. NDQ (Wang et al., 2020b) proposes nearly decomposable value functions to address the miscoordination problem in learning fully decentralized value functions. Emergence is a topic with increasing interest in deep MARL. Works on the emergence of communication (Foerster et al., 2016; Lazaridou et al., 2017; Das et al., 2017; Mordatch & Abbeel, 2018; Wang et al.,

2020b; Kang et al., 2020), the emergence of fairness (Jiang & Lu, 2019), and the emergence of tool usage (Baker et al., 2020) provide a deep learning perspective in understanding both natural and artificial multi-agent systems.

To learn diverse and identifiable roles, we propose to optimize the mutual information between individual roles and trajectories. A recent work studying multi-agent exploration, MAVEN (Mahajan et al., 2019), uses a similar objective. Different from ROMA, MAVEN aims at committed exploration. This difference in high-level purpose leads to many technical distinctions. First, MAVEN optimizes the mutual information between the joint trajectory and a latent variable conditioned on a Gaussian or uniform random variable to encourage diverse joint trajectory. Second, apart from the mutual information objective, we propose a novel regularizer to learn specialized roles, while MAVEN adopts a hierarchical structure and encourages the latent variable to help get more environmental rewards. We empirically compare ROMA with MAVEN in Sec. 5. More related works will be discussed in Appendix D.

5. Experiments

Our experiments aim to answer the following questions: (1) Whether the learned roles can automatically adapt in dynamic environments? (Sec. 5.1.) (2) Can our method promote sub-task specialization? That is, agents with similar responsibilities have similar role embedding representations, while agents with different responsibilities have role embedding representations far from each other. (Sec. 5.1, 5.3.) (3) Can such sub-task specialization improve the performance of multi-agent reinforcement learning algorithms? (Sec. 5.2.) (4) How do roles evolve during training, and how do they influence team performance? (Sec. 5.4.) (5) Can the dissimilarity model d_ϕ learn to measure the dissimilarity

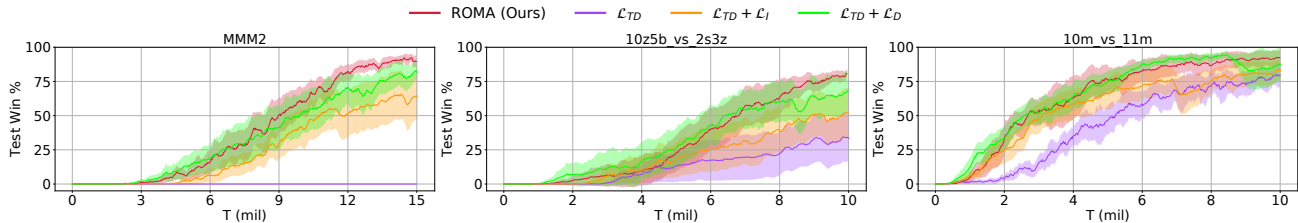


Figure 5. Ablation studies regarding the two role-learning losses.

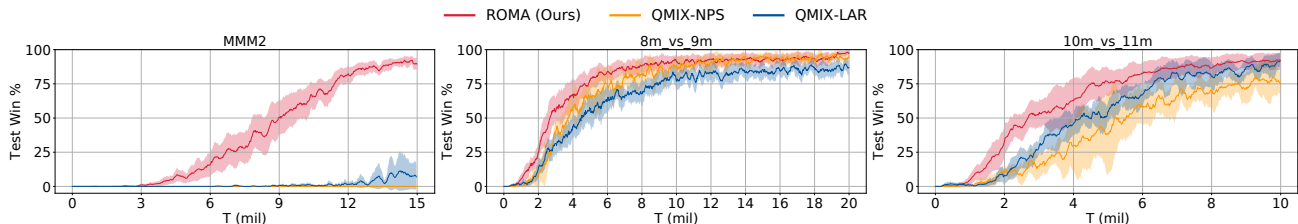


Figure 6. Comparison of our method against ablations.

between agents’ trajectories? (Sec. 5.4.) Videos² of our experiments and the code³ are available online.

Baselines We compare our methods with various baselines shown in Table 1. In particular, we carry out the following ablation studies: (i) We separately omit each (or both) of the two role-learning objectives (\mathcal{L}_I and \mathcal{L}_D) while leaving the other parts of ROMA unchanged. These three ablations are designed to highlight the contribution of each of the proposed regularizers. (ii) QMIX-NPS. The same as QMIX (Rashid et al., 2018), but agents do not share parameters. Our method achieves adaptive learning sharing, and comparison against QMIX (parameters are shared among agents) and QMIX-NPS tests whether this flexibility can improve learning efficiency. (iii) QMIX-LAR, QMIX with a similar number of parameters with our framework, which can test whether the superiority of our method comes from the increase in the number of parameters.

We carry out a grid search over the loss coefficients λ_I and λ_D , and fix them at 10^{-4} and 10^{-2} , respectively, across all the experiments. The dimensionality of latent role space is set to 3, so we did not use any dimensionality reduction techniques when visualizing the role embedding representations. Other hyperparameters are also fixed in our experiments, which are listed in Appendix B.1. For ROMA, We use elementary network structures (fully-connected networks or GRU) for the role encoder, role decoder, and trajectory encoder. The details of the architecture of our method and baselines can be found in Appendix B.

²<https://sites.google.com/view/romarl/>

³<https://github.com/TonghanWang/ROMA>

Table 1. Baseline algorithms.

| | Alg. | Description |
|---------------|------------------------------------|--|
| Related Works | IQL | Independent Q-learning |
| | COMA | Foerster et al. (2018) |
| | QMIX | Rashid et al. (2018) |
| | QTRAN | Son et al. (2019) |
| | MAVEN | Mahajan et al. (2019) |
| Ablations | \mathcal{L}_{TD} | ROMA without \mathcal{L}_I and \mathcal{L}_D |
| | $\mathcal{L}_{TD} + \mathcal{L}_I$ | ROMA without \mathcal{L}_D |
| | $\mathcal{L}_{TD} + \mathcal{L}_D$ | ROMA without \mathcal{L}_I |
| | QMIX-NPS | QMIX without parameter sharing among agents |
| | QMIX-LAR | QMIX with similar number of parameters with ROMA |

5.1. Dynamic Roles

Answering the first and second questions, we show snapshots in an episode played by ROMA agents on the StarCraft II micromanagement benchmark (SMAC) map 10m_vs_11m, where 10 Marines face 11 enemy Marines. As shown in Fig. 3 (the role representations at $t=27$ are presented in Fig. 1), although observations contain much information, such as positions, health points, shield points, states of ally and enemy units, etc., the role encoder learns to focus on different parts of the observations according to the dynamically changed situations. At the beginning ($t=1$), agents need to form a concave arc to maximize the number of agents whose shoot range covers the front line of enemies. ROMA learns to allocate roles according to agents’ relative positions so that agents can quickly form the offensive formation using

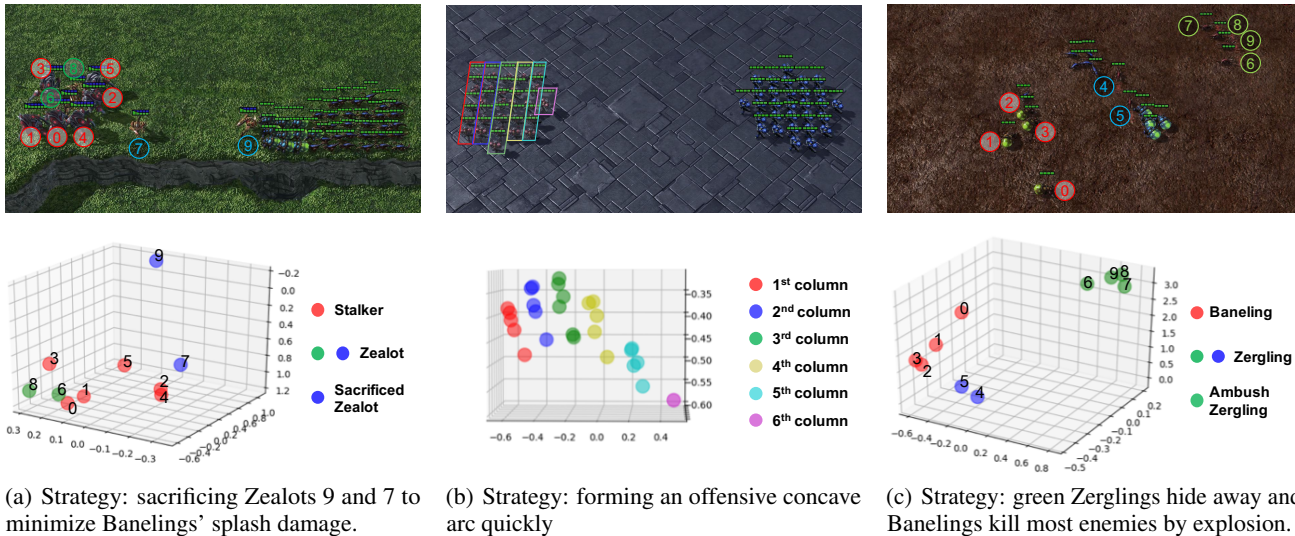


Figure 7. Learned roles for 6s4z_vs_10b30z, 27m_vs_30m, and 6z4b (means of the role distributions, μ_{ρ_i} , are shown, without using any dimensionality reduction techniques), and the related, automatically discovered responsibilities.

specialized policies. In the middle of the battle, one important tactic is to protect the injured ranged units. Our method learns this maneuver and roles cluster according to the remaining health points ($t=8, 19, 27$). Healthiest agents have role representations far from those of other agents. Such representations result in differentiated strategies: healthiest agents move forward to take on more firepower while other agents move backward, firing from a distance. In the meantime, some roles also cluster according to positions (agents 3 and 8 when $t=19$). The corresponding behaviors are agents with different roles fire alternatively to share the firepower. We can also observe that the role representations of dead agents aggregate together, representing a special group with an increasing number of agents during the battle.

These results demonstrate that our method learns dynamic roles and roles cluster clearly corresponding to automatically detected sub-tasks, in line with implicit constraints of the proposed optimization objectives.

5.2. Performance on StarCraft II

To test whether these roles and the corresponding sub-task specialization can improve learning efficiency, we test our method on the StarCraft II micromanagement (SMAC) benchmark (Samvelyan et al., 2019). This benchmark consists of various maps which have been classified as *easy*, *hard*, and *super hard*. We compare ROMA with algorithms shown in Table 1 and present results for one **easy** map (2s3z), three **hard** maps (5m_vs_6m, 8m_vs_9m & 10m_vs_11m), and two **super hard** maps (MMM2 & 27m_vs_30m). Although SMAC benchmark is challenging, it is not specially designed to test performance in tasks with

many agents. We thus introduce three new SMAC maps to test the scalability of our method, which are described in detail in Appendix C.

For evaluation, all experiments in this section are carried out with 5 different random seeds, and results are shown with a 95% confidence interval. Among these maps, four maps, MMM2, 6s4z_vs_10b30z, 6z4b, and 10z5b_vs_2z3s, feature heterogeneous agents, and the others have homogeneous agents. Fig. 4 shows that our method yields substantially better results than all the alternative approaches on both homogeneous and heterogeneous maps (additional plots can be found in Appendix C.1). MAVEN overcomes the negative effects of QMIX’s monotonicity constraint on exploration. However, it performs less satisfactorily than QMIX on most maps. We believe this is because agents start engaging in the battle immediately after spawning in SMAC maps, and exploration is not the critical factor affecting performance.

Ablations We carry out ablation studies, comparing with the ablations shown in Table 1 and present results on three maps: MMM2 (heterogeneous), 10z5b_vs_2s3z, and 10m_vs_11m (homogeneous) in Fig. 5 and 6. The superiority of our method against \mathcal{L}_{TD} highlights the contribution of the proposed regularizers – \mathcal{L}_{TD} performs even worse than QMIX on two of the three maps. By comparing ROMA with $\mathcal{L}_{TD} + \mathcal{L}_I$ and $\mathcal{L}_{TD} + \mathcal{L}_D$, we can conclude that the specialization loss \mathcal{L}_D is more important in terms of performance improvements. Introducing \mathcal{L}_I can make training more stable (for example, on the map 10m_vs_11m), but optimizing \mathcal{L}_I alone can only slightly improve the performance. These observations support the claim that sub-task specialization can improve labor efficiency.

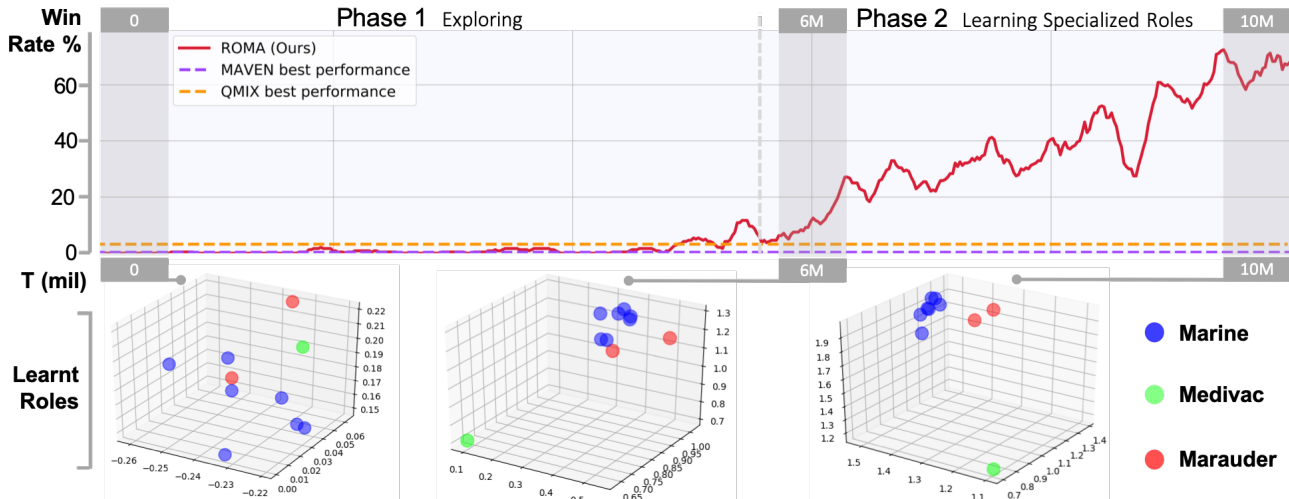


Figure 8. Role emergence and evolution on the map MMM2 (role representations at time step 1 are shown) during training (means of the role distributions, μ_{ρ_i} , are shown, without using any dimensionality reduction techniques). The emergence and specialization of roles is closely connected to the improvement of team performance. Agents in MMM2 are heterogeneous, and we show role evolution process in a homogeneous team in Appendix C.3.

Comparison between QMIX-NPS and QMIX demonstrates that parameter sharing can, as documented (Foerster et al., 2018; Rashid et al., 2018), speed up training. As discussed in the introduction, both these two paradigms may not get the best possible performance. In contrast, our method provides a dynamic learning sharing mechanism – agents committed to a certain responsibility have similar policies. The comparison of the performance of ROMA, QMIX, and QMIX-NPS proves that such sub-task specialization can indeed improve team performance. What’s more, comparison of ROMA against QMIX-LAR proves that the superiority of our method does not depend on the larger number of parameters.

The performance gap between ROMA and ablations is more significant on maps with more than ten agents. This observation supports discussions in previous sections – the emergence of role is more likely to improve the labor efficiency in larger populations.

5.3. Role Embedding Representations

To explain the superiority of ROMA, we present the learned role embedding representations for three maps in Fig. 7. Roles are representative of automatically discovered sub-tasks in the learned winning strategy. In the map of 6s4z_vs_10b30z, ROMA learns to sacrifice Zealots 9 and 7 to kill all the enemy Banelings. Specifically, Zealots 9 and 7 will move to the frontier one by one to minimize the splash damage, while other agents will stay away and wait until all Banelings explode. Fig. 7(a) shows the role embedding representations while performing the first sub-task where

agent 9 is sacrificed. We can see that the role of Zealot 9 is quite different from those of other agents. Correspondingly, the strategy at this time is agent 9 moving rightward while other agents keep still. Detailed analysis for the other two maps can be found in Appendix C.2.

5.4. Emergence and Evolution of Roles

We have shown the learned role representations and performance of our method, but the relationship between roles and performance remains unclear. To make up for this shortcoming, we visualize the emergence and evolution of roles during the training process on the map MMM2 (heterogeneous) and 10m_vs_11m (homogeneous). We discuss the results on MMM2 here and defer analysis of 10m_vs_11m to Appendix C.3.

In MMM2, 1 Medivac, 2 Marauders, and 7 Marines are faced with a stronger enemy team consisting of 1 Medivac, 3 Marauders, and 8 Marines. Among the three involved unit types, Medivac is the most special one for that it can heal the injured units. In Fig. 8, we show one of the learning curves of ROMA (red) and the role representations at the first environment step at three different stages. When the training begins ($T=0$), roles are random, and the agents are exploring the environment to learn the basic dynamics and the structure of the task. By $T=6M$, ROMA has learned that the responsibilities of the Medivac are different from those of Marines and Marauders. The role, and correspondingly, the policy of the Medivac becomes quite different (Fig. 8 middle). Such differentiation in behaviors enables agents to start winning the game. Gradually, ROMA learns that

Table 2. The mean and standard deviation of the learned dissimilarities d_ϕ between agents’ trajectories on the map MMM2.

| | |
|------------------------------|---------------------|
| Between different unit types | 0.9556 ± 0.0009 |
| Between the same unit type | 0.0780 ± 0.0019 |

Marines and Marauders have dissimilar characteristics and should take different sub-tasks, indicated by the differentiation of their role representations (Fig. 8 right). This further specialization facilitates the performance increase between 6M and 10M. After $T=10M$, the responsibilities of roles are clear, and, as a result, the win rate gradually converges (Fig. 4 top left). For comparison, ROMA without \mathcal{L}_I and \mathcal{L}_D can not even win once on this challenging task (\mathcal{L}_{TD} in Fig. 6-left). These results demonstrate that the gradually specialized roles are indispensable in team performance improvement.

Moreover, we find that the learned dissimilarity model d_ϕ introduced in Sec. 3.2 provides an empirical evaluation for identifying new roles. We use the map MMM2 as an example, where, as we discussed above, the learned roles of agents are characterized by their unit types. After scaling to $[0, 1]$, the learned dissimilarity between trajectories of agents with different unit types is close to 0.96, while the learned dissimilarity between trajectories of agents with the same unit type is around 0.08. These results indicate that an appropriate threshold can be used to decide when an individual behavior (trajectory) can be assigned the terminology *role*.

In summary, our experiments demonstrate that ROMA can learn dynamic, identifiable, versatile, and specialized roles that effectively decompose the task. Drawing support from these emergent roles, our method significantly pushes forward the state of the art of multi-agent reinforcement learning algorithms.

6. Closing Remarks

We have introduced the concept of roles into deep multi-agent reinforcement learning by capturing the emergent roles and encouraging them to specialize on a set of automatically detected sub-tasks. Such deep role-oriented multi-agent learning framework provides another perspective to explain and promote cooperation within agent teams, and implicitly draws connection to the division of labor, which has been practiced in many natural systems for long.

To our best knowledge, this paper is making a first attempt at learning roles via deep reinforcement learning. The gargantuan task of understanding the emergence of roles, the division of labor, and interactions between more complex roles in hierarchical organization still lies ahead. We believe that these topics are basic and indispensable in building effective, flexible, and general-purpose multi-agent systems

and this paper can help tackle these challenges.

Acknowledgements

We gratefully acknowledge Jin Zhang for his valuable discussions. We also would like to thank the reviewers for their detailed and constructive feedback. This work is partially supported by the sponsorship of Guoqiang Institute at Tsinghua University.

References

- Alemi, A. A., Fischer, I., Dillon, J. V., and Murphy, K. Deep variational information bottleneck. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- Baker, B., Kanitscheider, I., Markov, T., Wu, Y., Powell, G., McGrew, B., and Mordatch, I. Emergent tool use from multi-agent autocurricula. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- Bargiacchi, E., Verstraeten, T., Roijers, D., Nowé, A., and Hasselt, H. Learning to coordinate with coordination graphs in repeated single-stage multi-agent decision problems. In *International Conference on Machine Learning*, pp. 491–499, 2018.
- Becht, M., Gurzki, T., Klarmann, J., and Muscholl, M. Rope: Role oriented programming environment for multiagent systems. In *Proceedings Fourth IFCIS International Conference on Cooperative Information Systems. CoopIS 99 (Cat. No. PR00384)*, pp. 325–333. IEEE, 1999.
- Bonjean, N., Mefteh, W., Gleizes, M., Maurel, C., and Migeon, F. Adelfe 2.0: Handbook on agent-oriented design processes, m. cossentino, v. hilaire, a. molesini, and v. seidita, 2014.
- Bresciani, P., Perini, A., Giorgini, P., Giunchiglia, F., and Mylopoulos, J. Tropos: An agent-oriented software development methodology. *Autonomous Agents and Multi-Agent Systems*, 8(3):203–236, 2004.
- Butler, E. *The condensed wealth of nations*. Centre for Independent Studies, 2012.
- Cao, Y., Yu, W., Ren, W., and Chen, G. An overview of recent progress in the study of distributed multi-agent coordination. *IEEE Transactions on Industrial informatics*, 9(1):427–438, 2012.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *Proceedings of the 2014*

- Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734, 2014.
- Cossentino, M., Gaglio, S., Sabatucci, L., and Seidita, V. The passi and agile passi mas meta-models compared with a unifying proposal. In *International Central and Eastern European Conference on Multi-Agent Systems*, pp. 183–192. Springer, 2005.
- Das, A., Kottur, S., Moura, J. M., Lee, S., and Batra, D. Learning cooperative visual dialog agents with deep reinforcement learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2951–2960, 2017.
- Das, A., Gervet, T., Romoff, J., Batra, D., Parikh, D., Rabbat, M., and Pineau, J. Tarmac: Targeted multi-agent communication. In *International Conference on Machine Learning*, pp. 1538–1546, 2019.
- DeLoach, S. A. and Garcia-Ojeda, J. C. O-mase: a customisable approach to designing and building complex, adaptive multi-agent systems. *International Journal of Agent-Oriented Software Engineering*, 4(3):244–280, 2010.
- Depke, R., Heckel, R., and Küster, J. M. Roles in agent-oriented modeling. *International Journal of Software engineering and Knowledge engineering*, 11(03):281–302, 2001.
- Ferber, J., Gutknecht, O., and Michel, F. From agents to organizations: an organizational view of multi-agent systems. In *International workshop on agent-oriented software engineering*, pp. 214–230. Springer, 2003.
- Foerster, J., Assael, I. A., de Freitas, N., and Whiteson, S. Learning to communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 2137–2145, 2016.
- Foerster, J., Nardelli, N., Farquhar, G., Afouras, T., Torr, P. H., Kohli, P., and Whiteson, S. Stabilising experience replay for deep multi-agent reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1146–1155. JMLR. org, 2017.
- Foerster, J. N., Farquhar, G., Afouras, T., Nardelli, N., and Whiteson, S. Counterfactual multi-agent policy gradients. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Gordon, D. M. The organization of work in social insect colonies. *Nature*, 380(6570):121–124, 1996.
- Grover, A., Al-Shedivat, M., Gupta, J. K., Burda, Y., and Edwards, H. Evaluating generalization in multiagent systems using agent-interaction graphs. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 1944–1946. International Foundation for Autonomous Agents and Multiagent Systems, 2018.
- Guestrin, C., Koller, D., and Parr, R. Multiagent planning with factored mdps. In *Advances in neural information processing systems*, pp. 1523–1530, 2002.
- Hoshen, Y. Vain: Attentional multi-agent predictive modeling. In *Advances in Neural Information Processing Systems*, pp. 2701–2711, 2017.
- Iqbal, S. and Sha, F. Actor-attention-critic for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 2961–2970, 2019.
- Jeanson, R., Kukuk, P. F., and Fewell, J. H. Emergence of division of labour in halictine bees: contributions of social interactions and behavioural variance. *Animal behaviour*, 70(5):1183–1193, 2005.
- Jiang, J. and Lu, Z. Learning attentional communication for multi-agent cooperation. In *Advances in Neural Information Processing Systems*, pp. 7254–7264, 2018.
- Jiang, J. and Lu, Z. Learning fairness in multi-agent systems. In *Advances in Neural Information Processing Systems*, pp. 13854–13865, 2019.
- Kang, Y., Wang, T., and de Melo, G. Incorporating pragmatic reasoning communication into emergent language. *arXiv preprint arXiv:2006.04109*, 2020.
- Kim, D., Moon, S., Hostallero, D., Kang, W. J., Lee, T., Son, K., and Yi, Y. Learning to schedule communication in multi-agent reinforcement learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Kipf, T., Fetaya, E., Wang, K.-C., Welling, M., and Zemel, R. Neural relational inference for interacting systems. In *International Conference on Machine Learning*, pp. 2688–2697, 2018.
- Lazaridou, A., Peysakhovich, A., and Baroni, M. Multi-agent cooperation and the emergence of (natural) language. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- Lhaksmana, K. M., Murakami, Y., and Ishida, T. Role-based modeling for designing agent behavior in self-organizing multi-agent systems. *International Journal of Software Engineering and Knowledge Engineering*, 28(01):79–96, 2018.
- Lowe, R., Wu, Y., Tamar, A., Harb, J., Abbeel, O. P., and Mordatch, I. Multi-agent actor-critic for mixed

- cooperative-competitive environments. In *Advances in Neural Information Processing Systems*, pp. 6379–6390, 2017.
- Mahajan, A., Rashid, T., Samvelyan, M., and Whiteson, S. Maven: Multi-agent variational exploration. In *Advances in Neural Information Processing Systems*, pp. 7611–7622, 2019.
- Mordatch, I. and Abbeel, P. Emergence of grounded compositional language in multi-agent populations. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Nguyen, D. T., Kumar, A., and Lau, H. C. Credit assignment for collective multiagent rl with global rewards. In *Advances in Neural Information Processing Systems*, pp. 8102–8113, 2018.
- Nowé, A., Vrancx, P., and De Hauwere, Y.-M. Game theory and multi-agent reinforcement learning. In *Reinforcement Learning*, pp. 441–470. Springer, 2012.
- Odell, J., Nodine, M., and Levy, R. A metamodel for agents, roles, and groups. In *International Workshop on Agent-Oriented Software Engineering*, pp. 78–92. Springer, 2004.
- Oliehoek, F. A., Amato, C., et al. *A concise introduction to decentralized POMDPs*, volume 1. Springer, 2016.
- Omicini, A. Soda: Societies and infrastructures in the analysis and design of agent-based systems. In *International Workshop on Agent-Oriented Software Engineering*, pp. 185–193. Springer, 2000.
- Padgham, L. and Winikoff, M. Prometheus: A methodology for developing intelligent agents. In *International Workshop on Agent-Oriented Software Engineering*, pp. 174–185. Springer, 2002.
- Pavón, J. and Gómez-Sanz, J. Agent oriented software engineering with ingenias. In *International Central and Eastern European Conference on Multi-Agent Systems*, pp. 394–403. Springer, 2003.
- Peng, P., Wen, Y., Yang, Y., Yuan, Q., Tang, Z., Long, H., and Wang, J. Multiagent bidirectionally-coordinated nets: Emergence of human-level coordination in learning to play starcraft combat games. *arXiv preprint arXiv:1703.10069*, 2017.
- Rashid, T., Samvelyan, M., Witt, C. S., Farquhar, G., Foerster, J., and Whiteson, S. Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 4292–4301, 2018.
- Samvelyan, M., Rashid, T., de Witt, C. S., Farquhar, G., Nardelli, N., Rudner, T. G., Hung, C.-M., Torr, P. H., Foerster, J., and Whiteson, S. The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- Singh, A., Jain, T., and Sukhbaatar, S. Learning when to communicate at scale in multiagent cooperative and competitive tasks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Smith, A. The wealth of nations [1776], 1937.
- Son, K., Kim, D., Kang, W. J., Hostallero, D. E., and Yi, Y. Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *International Conference on Machine Learning*, pp. 5887–5896, 2019.
- Spanoudakis, N. and Moraitis, P. Using aseme methodology for model-driven agent systems development. In *International Workshop on Agent-Oriented Software Engineering*, pp. 106–127. Springer, 2010.
- Stone, P. and Veloso, M. Task decomposition, dynamic role assignment, and low-bandwidth communication for real-time strategic teamwork. *Artificial Intelligence*, 110(2):241–273, 1999.
- Sukhbaatar, S., Fergus, R., et al. Learning multiagent communication with backpropagation. In *Advances in Neural Information Processing Systems*, pp. 2244–2252, 2016.
- Sunehag, P., Lever, G., Gruslys, A., Czarniecki, W. M., Zambaldi, V., Jaderberg, M., Lanctot, M., Sonnerat, N., Leibo, J. Z., Tuyls, K., et al. 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*, pp. 2085–2087. International Foundation for Autonomous Agents and Multiagent Systems, 2018.
- Tan, M. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the tenth international conference on machine learning*, pp. 330–337, 1993.
- Usunier, N., Synnaeve, G., Lin, Z., and Chintala, S. Episodic exploration for deep deterministic policies: An application to starcraft micromanagement tasks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- Vinyals, O., Ewalds, T., Bartunov, S., Georgiev, P., Vezhnevets, A. S., Yeo, M., Makhzani, A., Küttler, H., Agapiou, J., Schrittwieser, J., et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.

- Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., Choi, D. H., Powell, R., Ewalds, T., Georgiev, P., et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575 (7782):350–354, 2019.
- Wainwright, M. J., Jordan, M. I., et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- Wang, T., Wang, J., Yi, W., and Zhang, C. Influence-based multi-agent exploration. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020a.
- Wang, T., Wang, J., Zheng, C., and Zhang, C. Learning nearly decomposable value functions with communication minimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020b.
- Wen, Y., Yang, Y., Luo, R., Wang, J., and Pan, W. Probabilistic recursive reasoning for multi-agent reinforcement learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- Wooldridge, M., Jennings, N. R., and Kinny, D. The gaia methodology for agent-oriented analysis and design. *Autonomous Agents and multi-agent systems*, 3(3):285–312, 2000.
- Yang, Z., Zhao, J., Dhingra, B., He, K., Cohen, W. W., Salakhutdinov, R. R., and LeCun, Y. Glomo: unsupervised learning of transferable relational graphs. In *Advances in Neural Information Processing Systems*, pp. 8950–8961, 2018.
- Zhang, C. and Lesser, V. Coordinated multi-agent reinforcement learning in networked distributed pomdps. In *Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.
- Zhu, H. and Zhou, M. Role-based multi-agent systems. In *Personalized Information Retrieval and Access: Concepts, Methods and Practices*, pp. 254–285. Igi Global, 2008.

A. Mathematical Derivation

A.1. Identifiable Roles

For learning identifiable roles, we propose to maximize the conditional mutual information objective between roles and local observation-action histories given the current observations. In Sec. 3.1 of the paper, we introduce a posterior estimator and derive a tractable lower bound of the mutual information term:

$$\begin{aligned}
 I(\rho_i^t; \tau_i^{t-1} | o_i^t) &= \mathbb{E}_{\rho_i^t, \tau_i^{t-1}, o_i^t} \left[\log \frac{p(\rho_i^t | \tau_i^{t-1}, o_i^t)}{p(\rho_i^t | o_i^t)} \right] \\
 &= \mathbb{E}_{\rho_i^t, \tau_i^{t-1}, o_i^t} \left[\log \frac{q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)}{p(\rho_i^t | o_i^t)} \right] \\
 &\quad + \mathbb{E}_{\tau_i^{t-1}, o_i^t} [D_{\text{KL}}(p(\rho_i^t | \tau_i^{t-1}, o_i^t) \| q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t))] \\
 &\geq \mathbb{E}_{\rho_i^t, \tau_i^{t-1}, o_i^t} \left[\log \frac{q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)}{p(\rho_i^t | o_i^t)} \right],
 \end{aligned} \tag{8}$$

where the last inequality holds because of the non-negativity of the KL divergence. Then it follows that:

$$\begin{aligned}
 &\mathbb{E}_{\rho_i^t, \tau_i^{t-1}, o_i^t} \left[\log \frac{q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)}{p(\rho_i^t | o_i^t)} \right] \\
 &= \mathbb{E}_{\rho_i^t, \tau_i^{t-1}, o_i^t} [\log q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)] - \mathbb{E}_{\rho_i^t, o_i^t} [\log p(\rho_i^t | o_i^t)] \\
 &= \mathbb{E}_{\rho_i^t, \tau_i^{t-1}, o_i^t} [\log q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)] + \mathbb{E}_{o_i^t} [H(\rho_i^t | o_i^t)] \\
 &= \mathbb{E}_{\tau_i^{t-1}, o_i^t} \left[\int p(\rho_i^t | \tau_i^{t-1}, o_i^t) \log q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t) d\rho_i^t \right] + \mathbb{E}_{o_i^t} [H(\rho_i^t | o_i^t)]
 \end{aligned} \tag{9}$$

The role encoder is conditioned on the local observations, so given the observations, the distributions of roles, $p(\rho_i^t)$, are independent from the local histories. Thus, we have

$$I(\rho_i^t; \tau_i^{t-1} | o_i^t) \geq -\mathbb{E}_{\tau_i^{t-1}, o_i^t} [\mathcal{CE}[p(\rho_i^t | o_i^t) \| q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)]] + \mathbb{E}_{o_i^t} [H(\rho_i^t | o_i^t)] \tag{10}$$

In practice, we use a replay buffer \mathcal{D} and minimize

$$\mathcal{L}_I(\theta_\rho, \xi) = \mathbb{E}_{(\tau_i^{t-1}, o_i^t) \sim \mathcal{D}} [\mathcal{CE}[p(\rho_i^t | o_i^t) \| q_\xi(\rho_i^t | \tau_i^{t-1}, o_i^t)] - H(\rho_i^t | o_i^t)]. \tag{11}$$

A.2. Specialized Roles

Conditioning roles on local observations enables roles to be dynamic, and optimizing \mathcal{L}_I enables roles to be identifiable by agents' long-term behaviors, but these formulations do not explicitly encourage specialized roles. To make up for this shortcoming, we propose a role differentiation objective in Sec. 3.2 of the paper, where a mutual information maximization objective is involved (maximizing $I(\rho_i^t; \tau_j^{t-1} | o_j^t)$). Here, we derive a variational lower bound of this mutual information objective to render it feasible to be optimized.

$$\begin{aligned}
 I(\rho_i^t; \tau_j^{t-1} | o_j^t) &= \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} \left[\log \frac{p(\rho_i^t, \tau_j^{t-1} | o_j^t)}{p(\rho_i^t | o_j^t) p(\tau_j^{t-1} | o_j^t)} \right] \\
 &= \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} \left[\log \frac{p(\rho_i^t | \tau_j^{t-1}, o_j^t)}{p(\rho_i^t | o_j^t)} \right] \\
 &= \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} [\log p(\rho_i^t | \tau_j^{t-1}, o_j^t)] + \mathbb{E}_{o_j^t} [H(\rho_i^t | o_j^t)] \\
 &\geq \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} [\log p(\rho_i^t | \tau_j^{t-1}, o_j^t)].
 \end{aligned} \tag{12}$$

We clip the variances of role distributions at a small value (0.1) to ensure that the entropy of role distributions are always

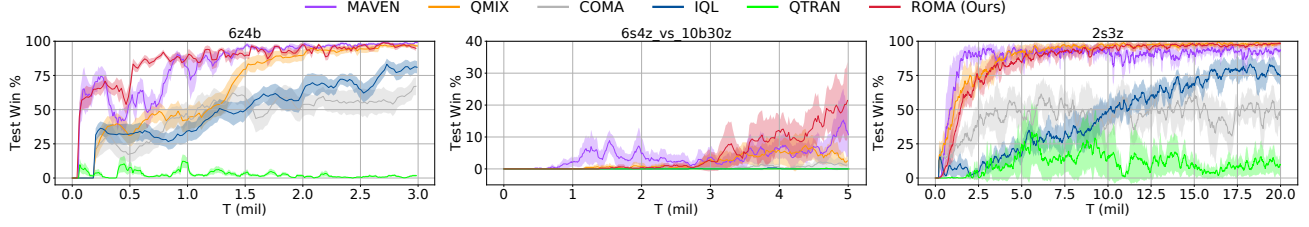


Figure 9. Additional results on the SMAC benchmark.

non-negative so that the last inequality holds. Then, it follows that:

$$\begin{aligned}
 & \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} [\log p(\rho_i^t | \tau_j^{t-1}, o_j^t)] \\
 &= \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} [\log q_\xi(\rho_i^t | \tau_j^{t-1}, o_j^t)] + \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} [D_{\text{KL}} [p(\rho_i^t | \tau_j^{t-1}, o_j^t) \| q_\xi(\rho_i^t | \tau_j^{t-1}, o_j^t)]] \\
 &\geq \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} [\log q_\xi(\rho_i^t | \tau_j^{t-1}, o_j^t)],
 \end{aligned} \tag{13}$$

where q_ξ is the trajectory encoder introduced in Sec. A.1, and the KL divergence term can be left out when deriving the lower bound because it is non-negative. Therefore, we have:

$$I(\rho_i^t; \tau_j^{t-1} | o_j^t) \geq \mathbb{E}_{\rho_i^t, \tau_j^{t-1}, o_j^t} [\log q_\xi(\rho_i^t | \tau_j^{t-1}, o_j^t)]. \tag{14}$$

Recall that, in order to learn specialized roles, we propose to minimize:

$$\|D_\phi^t\|_F - \sum_{i \neq j} \min\{I(\rho_i^t; \tau_j^{t-1} | o_j^t) + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\}, \tag{15}$$

where $D_\phi^t = (d_{ij}^t)$, and $d_{ij}^t = d_\phi(\tau_i^{t-1}, \tau_j^{t-1})$ is the estimated dissimilarity between trajectories of agent i and j . For the term $\min\{I(\rho_i^t; \tau_j^{t-1} | o_j^t) + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\}$, we have:

$$\begin{aligned}
 & \min\{I(\rho_i^t; \tau_j^{t-1} | o_j^t) + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\} \\
 &= \min\{\mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} \left[\log \frac{p(\rho_i^t, \tau_j^{t-1} | o_j^t)}{p(\rho_i^t | o_j^t) p(\tau_j^{t-1} | o_j^t)} + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}) \right], \mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [U]\},
 \end{aligned} \tag{16}$$

where τ^{t-1} is the joint trajectory, \mathbf{o}^t is the joint observation, and $\boldsymbol{\rho}^t = \langle \rho_1^t, \rho_2^t, \dots, \rho_n^t \rangle$. We denote

$$\begin{aligned}
 T_1 &\equiv \log \frac{p(\rho_i^t, \tau_j^{t-1} | o_j^t)}{p(\rho_i^t | o_j^t) p(\tau_j^{t-1} | o_j^t)}, \\
 T_2 &\equiv \log q_\xi(\rho_i^t | \tau_j^{t-1}, o_j^t).
 \end{aligned} \tag{17}$$

Because

$$\begin{aligned}
 T_2 &\geq \min\{T_2, U\}, \\
 U &\geq \min\{T_2, U\},
 \end{aligned} \tag{18}$$

it follows that:

$$\begin{aligned}
 \mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [T_2] &\geq \mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [\min\{T_2, U\}], \\
 \mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [U] &\geq \mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [\min\{T_2, U\}].
 \end{aligned} \tag{19}$$

So that

$$\begin{aligned}
 & \min\{\mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [T_1], \mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [U]\} \\
 &\geq \min\{\mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [T_2], \mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [U]\} \quad \{\text{Eq. 14}\} \\
 &\geq \mathbb{E}_{\tau^{t-1}, \mathbf{o}^t, \boldsymbol{\rho}^t} [\min\{T_2, U\}] \quad \{\text{Eq. 19}\},
 \end{aligned} \tag{20}$$

which means that Eq. 15 satisfies:

$$\|D_\phi^t\|_F - \sum_{i \neq j} \min\{I(\rho_i^t; \tau_j^{t-1} | o_j^t) + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\}$$

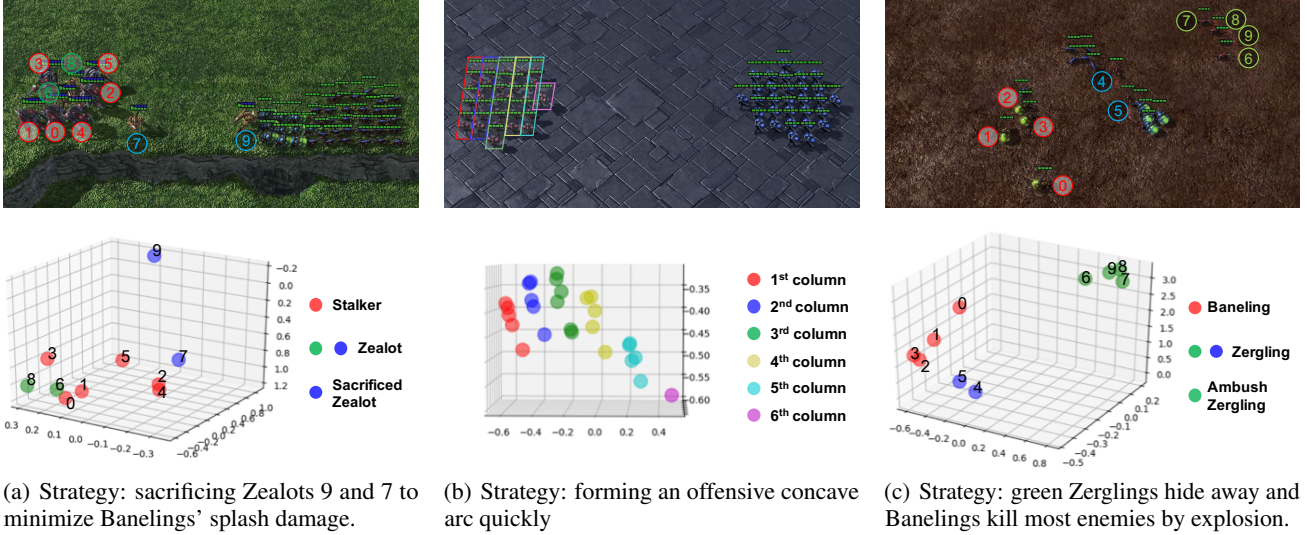


Figure 10. (Reproduced from Fig. 6 in the paper, for quick reference.) Learned roles for 6s4z_vs_10b30z, 27m_vs_30m, and 6z4b (means of the role distributions, μ_{ρ_i} , are shown, without using any dimensionality reduction techniques), and the related, automatically discovered responsibilities.

$$\begin{aligned}
 &= \mathbb{E}_{\tau^{t-1}, \mathcal{o}^t, \rho^t} [\|D_\phi^t\|_F] - \sum_{i \neq j} \min\{\mathbb{E}_{\tau^{t-1}, \mathcal{o}^t, \rho^t} [T_1 + d_\phi(\tau_i^{t-1}, \tau_j^{t-1})], U\} \\
 &\leq \mathbb{E}_{\tau^{t-1}, \mathcal{o}^t, \rho^t} [\|D_\phi^t\|_F] - \sum_{i \neq j} \mathbb{E}_{\tau^{t-1}, \mathcal{o}^t, \rho^t} [\min\{T_2 + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\}] \quad \{\text{Eq. 20}\} \\
 &= \mathbb{E}_{\tau^{t-1}, \mathcal{o}^t, \rho^t} \left[\left\| D_\phi^t \right\|_F - \sum_{i \neq j} \min\{T_2 + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\} \right].
 \end{aligned} \tag{21}$$

We minimize this upper bound to optimize Eq. 15. In practice, we use a replay buffer, and minimize:

$$\mathcal{L}_D(\theta_\rho, \phi, \xi) = \mathbb{E}_{(\tau^{t-1}, \mathcal{o}^t) \sim \mathcal{D}, \rho^t \sim p(\rho^t | \mathcal{o}^t)} \left[\left\| D_\phi^t \right\|_F - \sum_{i \neq j} \min\{q_\xi(\rho_i^t | \tau_j^{t-1}, \mathcal{o}_j^t) + d_\phi(\tau_i^{t-1}, \tau_j^{t-1}), U\} \right], \tag{22}$$

where \mathcal{D} is the replay buffer, τ^{t-1} is the joint trajectory, \mathcal{o}^t is the joint observation, and $\rho^t = \langle \rho_1^t, \rho_2^t, \dots, \rho_n^t \rangle$.

B. Architecture, Hyperparameters, and Infrastructure

B.1. ROMA

In this paper, we base our algorithm on QMIX (Rashid et al., 2018), whose framework is shown in Fig. 13 and described in Appendix D. In ROMA, each agent has a neural network to approximate its local utility. The local utility network consists of three layers, a fully-connected layer, followed by a 64 bit GRU, and followed by another fully-connected layer that outputs an estimated value for each action. The local utilities are fed into a mixing network estimating the global action value. The mixing network has a 32-dimensional hidden layer with ReLU activation. Parameters of the mixing network are generated by a hyper-net conditioning on the global state. This hyper-net has a fully-connected hidden layer of 32 dimensions. These settings are the same as QMIX.

We use very simple network structures for the components related to role embedding learning, i.e., the role encoder, the role decoder, and the trajectory encoder. The multi-variate Gaussian distributions from which the individual roles are drawn have their means and variances generated by the role encoder, which is a fully-connected network with a 12-dimensional hidden layer with ReLU activation. The parameters in the second fully-connected layers of the local utility approximators are generated by the role decoder whose inputs are the individual roles, which are 3-dimensional in all experiments. The role decoder is also a fully-connected network with a 12-dimensional hidden layer and ReLU activation. For the trajectory

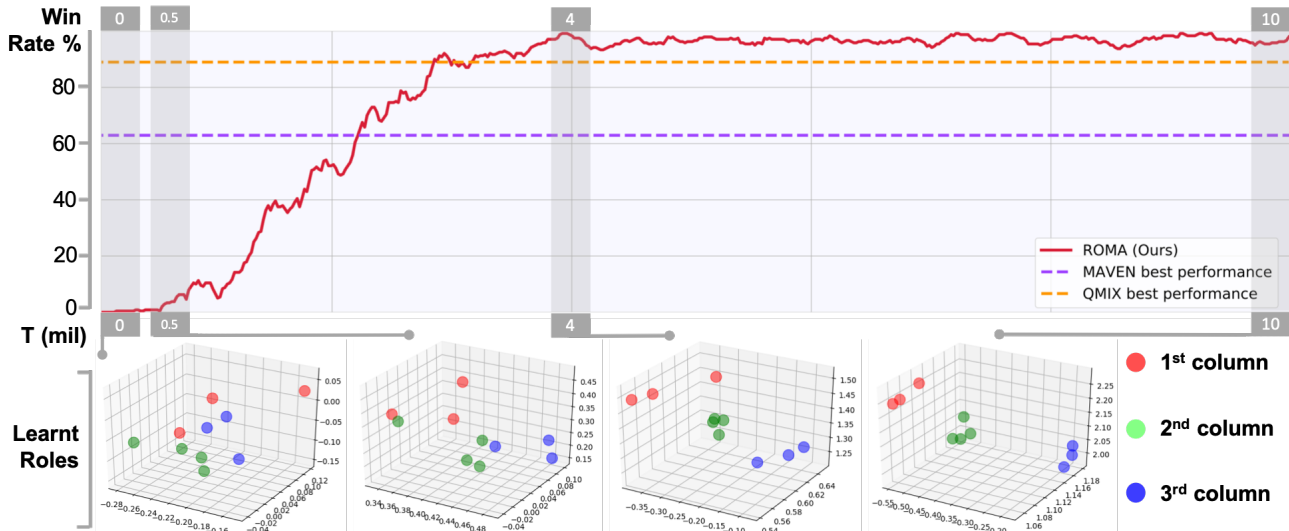


Figure 11. The process of role emergence and evolution on the map 10m_vs_11m.

encoder, we again use a fully-connected network with a 12-dimensional hidden layer and ReLU activation. The inputs of the trajectory encoder are the hidden states of the GRUs in the local utility functions after the last time step.

For all experiments, we set $\lambda_I = 10^{-4}$, $\lambda_D = 10^{-2}$, and the discount factor $\gamma = 0.99$. The optimization is conducted using RMSprop with a learning rate of 5×10^{-4} , α of 0.99, and with no momentum or weight decay. For exploration, we use ϵ -greedy with ϵ annealed linearly from 1.0 to 0.05 over $50k$ time steps and kept constant for the rest of the training. We run 8 parallel environments to collect samples. Batches of 32 episodes are sampled from the replay buffer, and the whole framework is trained end-to-end on fully unrolled episodes. All experiments on StarCraft II use the default reward and observation settings of the SMAC benchmark.

Experiments are carried out on NVIDIA GTX 2080 Ti GPU.

B.2. Baselines and Ablations

We compare ROMA with various baselines and ablations, which are listed in Table. 1 of the paper. For COMA (Foerster et al., 2018), QMIX (Rashid et al., 2018), and MAVEN (Mahajan et al., 2019), we use the codes provided by the authors where the hyper-parameters have been fin-tuned on the SMAC benchmark. QMIX-NPS uses the identical architecture as QMIX, and the only difference lies in that QMIX-NPS does not share parameters among agents. Compared to QMIX, for the local utility function of agents, QMIX-LAR adds two more fully-connected layers of 80 and 25 dimensions after the GRU layer so that it approximately has the same number of parameters as ROMA.

C. Additional Experimental Results

We benchmark our method on the StarCraft II unit micromanagement tasks. To test the scalability of the proposed approach, we introduce three maps. The 6z4b map features symmetry teams consisting of 4 Banelings and 6 Zerglings. In the map of 6s4z_vs_10b30z, 6 Stalkers and 4 Zealots learn to defeat 10 Banelings and 30 Zerglings. And 10z5b_vs_2z3s characterizes asymmetry teams consisting of 10 Zerglings & 5 Banelings and 2 Zealots & 3 Stalkers, respectively.

C.1. Performance Comparison against Baselines

Fig. 9 presents performance of ROMA against various baselines on three maps. Performance comparison on the other maps is shown in Fig. 4 of the paper. We can see that the advantage of ROMA is more significant on maps with more agents, such as 10z5b_vs_2z3s, MMM2, 27_vs_30m, and 10m_vs_11m.

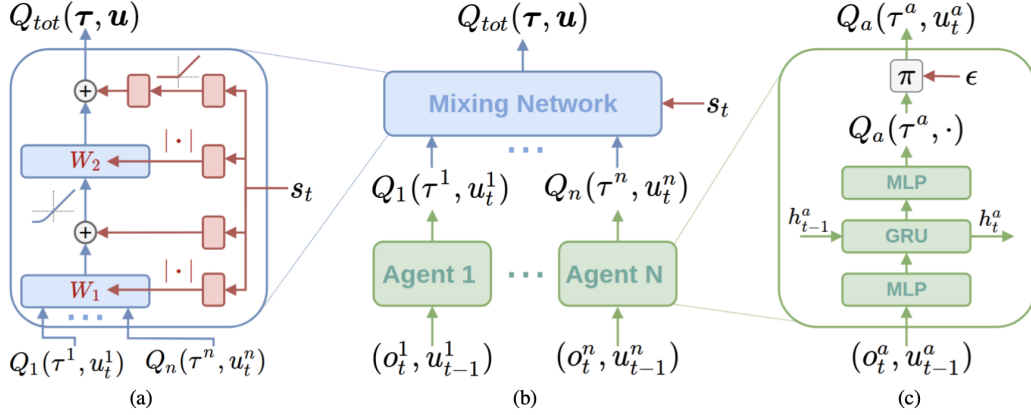


Figure 13. The framework of QMIX, reproduced from the original paper (Rashid et al., 2018). (a) The architecture of the mixing network (blue), whose weights and biases are generated by a hyper-net (red) conditioned on the global state. (b) The overall QMIX structure. (c) Local utility network structure.

C.2. Role Embedding Representations

Fig. 7 shows various roles learned by ROMA. Roles are closely related to the sub-tasks in the learned winning strategy.

For the map 27m_vs_30m, the winning strategy is to form an offensive concave arc before engaging in the battle. Fig. 10(b) illustrates the role embedding representations at the first time step when the agents are going to set up the attack formation. We can see the roles aggregate according to the relative positions of the agents. Such role differentiation leads to different moving strategies so that agents can quickly form the arc without collisions.

Similar role-behavior relationships can be seen in all tasks. We present another example on the task of 6z4b. In the winning strategy learned by ROMA, Zerglings 4 & 5 and Banelings kill most of the enemies, taking advantage of the splash damage of the Banelings, while Zerglings 6-9 hideaway, wait until the explosion is over, and then kill the remaining enemies. Fig. 10(c) shows the role embedding representations before the explosion. We can see clear clusters closely corresponding to the automatically detected sub-tasks at this time step.

Supported by these results, we can conclude that ROMA can automatically decompose the task and learn versatile roles, each of which is specialized in a certain sub-task.

C.3. Additional Results for Role Evolution

In Fig. 7 of the paper, we show how roles emerge and evolve on the map MMM2, where the involved agents are heterogeneous. In this section, we discuss the case of homogeneous agent teams. To this end, we visualize the emergence and evolution process of roles on the map 10m_vs_11m, which features 10 ally Marines facing 11 enemy Marines. In Fig. 11, we show the roles at the first time step of the battle (screenshot can be found in Fig. 12) at four different stages during the training. At this moment, agents need to form an offensive concave arc quickly. We can see that ROMA gradually learns to allocate roles according to relative positions of agents. Such roles and the corresponding differentiation in the individual policies help agents form the offensive arc more efficiently. Since setting up an attack formation is critical for winning the game, a connection between the specialization of the roles at the first time step and the improvement of the win rate can be observed.



Figure 12. Screenshot of 10m_vs_11m, $t=1$.

D. Related Works

Multi-agent reinforcement learning holds the promise to solve many real-world problems and has been making vigorous progress recently. To avoid otherwise exponentially large state-action space, factorizing MDPs for multi-agent systems is proposed (Guestrin et al., 2002). Coordination graphs (Bargiacchi et al., 2018; Yang et al., 2018; Grover et al., 2018; Kipf et al., 2018) and explicit communication (Sukhbaatar et al., 2016; Hoshen, 2017; Jiang & Lu, 2018; Singh et al., 2019; Das et al., 2019; Singh et al., 2019; Kim et al., 2019) are studied to model the dependence between the decision-making processes of agents. Training decentralized policies is faced with two challenges: the issue of non-stationarity (Tan, 1993) and reward assignment (Foerster et al., 2018; Nguyen et al., 2018). To resolve these problems, Sunehag et al. (2018) propose a value decomposition method called VDN. VDN learns a global action-value function, which is factored as the sum of each agent’s local Q-value. QMIX (Rashid et al., 2018) extends VDN by representing the global value function as a learnable, state-condition, and monotonic combination of the local Q-values. In this paper, we use the mixing network of QMIX. The framework of QMIX is shown in Fig. 13.

The StarCraft II unit micromanagement task is considered as one of the most challenging cooperative multi-agent testbeds for its high degree of control complexity and environmental stochasticity. Usunier et al. (2017) and Peng et al. (2017) study this problem from a centralized perspective. In order to facilitate decentralized control, we test our method on the SMAC benchmark (Samvelyan et al., 2019), which is the same as in (Foerster et al., 2017; 2018; Rashid et al., 2018; Mahajan et al., 2019).