# REVISITING THE MASTER-SLAVE ARCHITECTURE IN MULTI-AGENT DEEP REINFORCEMENT LEARNING

**Xiangyu Kong, Fangchen Liu**[∗] **& Yizhou Wang**
Nat'l Eng. Lab. for Video Technology, Cooperative Medianet Innovation Center,
Key Laboratory of Machine Perception (MoE), Sch'l of EECS,
Peking University
Beijing, 100871, China
`{kong,liufangchen,yizhou.wang}@pku.edu.cn`

**Bo Xin**[∗]
Microsoft Research
Beijing, 100080, China
`boxin@microsoft.com`

## ABSTRACT

Many tasks in artificial intelligence require the collaboration of multiple agents. We exam deep reinforcement learning for multi-agent domains. Recent research efforts often take the form of two seemingly conflicting perspectives, the decentralized perspective, where each agent is supposed to have its own controller; and the centralized perspective, where one assumes there is a larger model controlling all agents. In this regard, we revisit the idea of the master-slave architecture by incorporating both perspectives within one framework. Such a hierarchical structure naturally leverages advantages from one another. The idea of combining both perspectives is intuitive and can be well motivated from many real world systems, however, out of a variety of possible realizations, we highlights three key ingredients, i.e. composed action representation, learnable communication and independent reasoning. With network designs to facilitate these explicitly, our proposal consistently outperforms latest competing methods both in synthetic experiments and when applied to challenging StarCraft[1] micromanagement tasks.

## 1 INTRODUCTION

Reinforcement learning (RL) provides a formal framework concerned with how an agent takes actions in one environment so as to maximize some notion of cumulative reward. Recent years have witnessed successful application of RL technologies to many challenging problems, ranging from game playing Mnih et al. (2015); Silver et al. (2016) to robotics Levine et al. (2016) and other important artificial intelligence (AI) related fields such as Ren et al. (2017) etc. Most of these works have been studying the problem of a single agent.

However, many important tasks require the collaboration of multiple agents, for example, the coordination of autonomous vehicles Cao et al. (2013), multi-robot control Matignon et al. (2012), network packet delivery Ye et al. (2015) and multi-player games Synnaeve et al. (2016) to name a few. Although multi-agent reinforcement learning (MARL) methods have historically been applied in many settings Busoniu et al. (2008); Yang & Gu (2004), they were often restricted to simple environments and tabular methods.

Motivated from the success of (single agent) deep RL, where value/policy approximators were implemented via deep neural networks, recent research efforts on MARL also embrace deep networks and target at more complicated environments and complex tasks, e.g. Sukhbaatar et al. (2016); Peng et al. (2017); Foerster et al. (2017a); Mao et al. (2017) etc. Regardless though, it remains an open

---

[∗]Equal contribution.
[1]StarCraft and its expansion StarCraft: Brood War are trademarks of Blizzard Entertainment[TM]

challenge how deep RL can be effectively scaled to more agents in various situations. Deep RL is notoriously difficult to train. Moreover, the essential state-action space of multiple agents becomes geometrically large, which further exacerbates the difficulty of training for multi-agent deep reinforcement learning (deep MARL for short).

From the viewpoint of multi-agent systems, recent methods often take the form of one of two perspectives. That is, the decentralized perspective where each agent has its own controller; and the centralized perspective where there exists a larger model controlling all agents. As a consequence, learning can be challenging in the decentralized settings due to local viewpoints of agents, which perceive non-stationary environment due to concurrently exploring teammates. On the other hand, under a centralized perspective, one needs to directly deal with parameter search within the geometrically large state-action space originated from the combination of multiple agents.

In this regard, we revisit the idea of master-slave architecture to combine both perspectives in a complementary manner. The master-slave architecture is a canonical communication architecture which often effectively breaks down the original challenges of multiple agents. Such architectures have been well explored in multi-agent tasks Park et al. (2001); Verbeeck et al. (2005); Megherbi & Madera (2010); Megherbi & Kim (2016). Although our designs vary from these works, we have inherited the spirit of leveraging agent hierarchy in a master-slave manner. That is, the master agent tends to plan in a global manner without focusing on potentially distracting details from each slave agent and meanwhile the slave agents often locally optimize their actions with respect to both their local state and the guidance coming from the master agent. Such idea can be well motivated from many real world systems. One can consider the master agent as the central control of some organized traffic systems and the slave agents as each actual vehicles. Another instantiation of this idea is to consider the coach and the players in a football/basketball team. However, although the idea is clear and intuitive, we notice that our work is among the first to explicitly design master-slave architecture for deep MARL.

Specifically, we instantiate our idea with policy-based RL methods and propose a multi-agent policy network constructed with the master-slave agent hierarchy. For both each slave agent and the master agent, the policy approximators are realized using recurrent neural networks (RNN). At each time step, we can view the hidden states/representations of the recurrent cells as the "thoughts" of the agents. Therefore each agent has its own thinking/reasoning of the situation. While each slave agent takes local states as its input, the master agent takes both the global states and the messages from all slave agents as its input. The final action output of each slave agent is composed of contributions from both the corresponding slave agent and the master agent. This is implemented via a gated composition module (GCM) to process and transform "thoughts" from both agents to the final action.

We test our proposal (named MS-MARL) using both synthetic experiments and challenging Star-Craft micromanagement tasks. Our method consistently outperforms recent competing MARL methods by a clear margin. We also provide analysis to showcase the effectiveness of the learned policies, many of which illustrate interesting phenomena related to our specific designs.

In the rest of this paper, we first discuss some related works in Section 2. In Section 3, we introduce the detailed proposals to realize our master-slave multi-agent RL solution. Next, we move on to demonstrate the effectiveness of our proposal using challenging synthetic and real multi-agent tasks in Section 4. And finally Section 5 concludes this paper with discussions on our findings. Before proceeding, we summarize our major contributions as follows

- We revisit the idea of master-slave architecture for deep MARL. The proposed instantiation effectively combines both the centralized and decentralized perspectives of MARL.

- Our observations highlight and verify that composable action representation, independent master/slave reasoning and learnable communication in-between are key factors to be successful in MS-MARL.

- Our proposal empirically outperforms recent state-of-the-art methods on both synthetic experiments and challenging StarCraft micromanagement tasks, rendering it a novel competitive MARL solution in general.
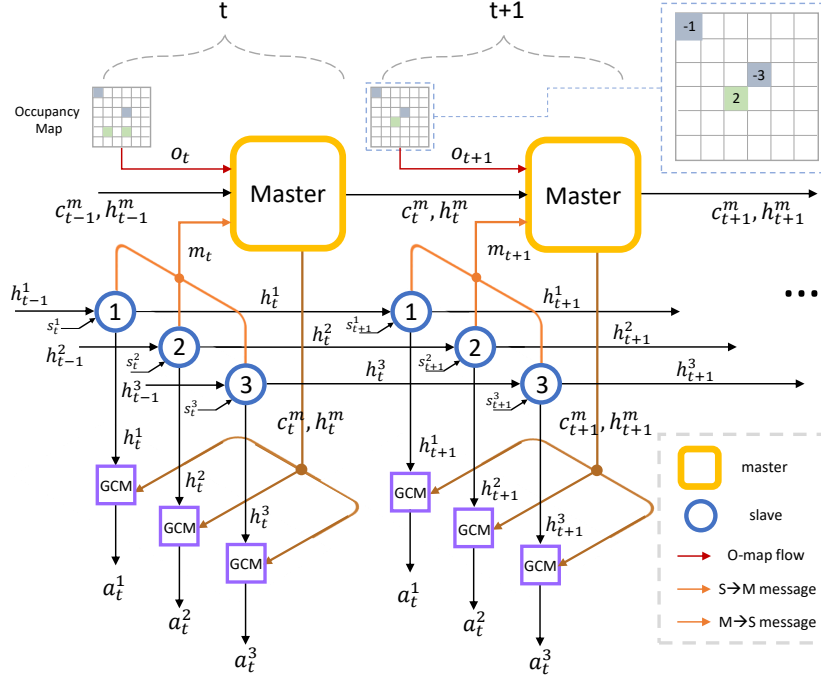
Figure 1: Pipeline of our master-slave multi-agent architecture

## 2 RELATED WORK

Current main stream RL methods apply conventional wisdoms such as Q-learning, policy gradient, actor-critic etc. Sutton & Barto (1998). Recent progress mainly focuses on practical adaptations especially when applying deep neural networks as value/policy approximators. Li (2017) provides a recent review on deep RL.

Although MARL has been studied in the past, they have been focused on simple tasks Busoniu et al. (2008). Only until recently, with the encouragement from the successes of deep RL, deep MARL has become a popular research area targeting at more complex and realistic tasks, see e.g. Foerster et al. (2016); Sukhbaatar et al. (2016); Kong et al. (2017); Peng et al. (2017); Foerster et al. (2017a); Mao et al. (2017) etc.

Foerster et al. (2016) and Kong et al. (2017) are among the first to propose learnable communications via back-propagation in deep Q-networks. However, due to their motivating tasks, both works focused on a decentralized perspective and usually applies to only a limited number of agents.

Usunier et al. (2016), Foerster et al. (2017b) and Peng et al. (2017) all proposed practical network structure or training strategies from a centralized perspective of MARL. Specifically, Peng et al. (2017) proposed a bidirectional communication channel among all agents to facilitate effective communication and many interesting designs toward the StarCraft micromanagement tasks. Usunier et al. (2016) proposed episodic exploration strategy for deterministic policy search and Foerster et al. (2017b) proposed the concept of stabilizing experience replay for MARL.

Note that the above works take only one of the two perspectives and are then inherently missing out the advantages of the other. Perhaps the most related works are from Foerster et al. (2016), Foerster et al. (2017a) and Mao et al. (2017). Sukhbaatar et al. (2016) proposed the "CommNet", where a broadcasting communication channel among all agents was set up to share global information realized as summation of the output from all individual agents. This design represents an initial version of the proposed master-slave framework, however the summed global signal is hand-crafted

3

information and moreover, this design does not facilitate an independently reasoning master agent. In Foerster et al. (2017a) and Mao et al. (2017), a global critic was proposed, which could potentially work at a centralized level, however since critics are basically value networks, they do not provide explicit policy guidance. Therefore they tend to work more like a commentator of a game who job is to analyze and criticize the play, rather than a coach coaching the game.

As discussed above, the master-slave architecture has already been studied in several multi-agent scenarios. Park et al. (2001) utilized the master-slave architecture to resolve conflicts between multiple soccer agents; while Verbeeck et al. (2005); Megherbi & Madera (2010); Megherbi & Kim (2016) explored master-slave hierarchy in RL applied to load-balancing and distributed computing environments. Our proposal can be viewed as a revisit to similar ideas for deep MARL. In this regard, we are among the first to combine both the centralized perspective and the decentralized perspective in an explicit manner. With the proposed designs, we facilitate independent master reasoning at a global level and each slave agent thinking at a local but focused scale, and collectively achieve optimal rewards via effective communication learned with back propagation. Compared with existing works, we emphasize such independent reasoning, the importance of which are well justified empirically in the experiments. We consistently outperform existing MARL methods and achieve state-of-the-art performance on challenging synthetic and real multi-agent tasks.

Since the master-slave architecture constructs agent hierarchy by definition, another interesting related field is hierarchical RL, e.g. Kulkarni et al. (2016); Vezhnevets et al. (2017). However, such hierarchical deep RL methods studies the hierarchy regarding tasks or goals and are usually targeting at sequential sub-tasks where the meta-controller constantly generates goals for controllers to achieve. Master-slave architecture, on the other hand, builds up hierarchy of multiple agents and mainly focuses on parallel agent-specific tasks instead, which is fundamentally different from the problems that hierarchical RL methods are concerned with.

## 3 MASTER-SLAVE MULTI-AGENT RL

We start by reiterating that the key idea is to facilitate both an explicit master controller that takes the centralized perspective and organize agents in a global or high level manner and all actual slave controllers work as the decentralized agents and optimize their specific actions relatively locally while depending on information from the master controller. Such an idea can be realized using either value-based methods, policy-based methods or actor-critic methods.

### 3.1 NETWORK ARCHITECTURE

Hereafter we focus on introducing an instantiation with policy gradient methods as an example, which also represents the actual solution in all our experiments. In particular, our target is to learn a mapping from states to actions $\pi_\theta(\boldsymbol{a}_t|\boldsymbol{s}_t)$ at any given time $t$, where $\boldsymbol{s} = \{s^m, s^1, ..., s^C\}$ and $\boldsymbol{a} = \{a^1, ..., a^C\}$ are collective states and actions of $C$ agents respectively and $\theta = \{\theta^m, \theta^1, ..., \theta^C\}$ represents the parameters of the policy function approximator of all agents, including the master agent $\theta^m$. Note that we have explicitly formulated $s^m$ to represent the independent state to the master agent but have left out a corresponding $a^m$ since the master's action will be merged with and represented by the final actions of all slave agents. This design has two benefits: 1) one can now input independent and potentially more global states to the master agent; and meanwhile 2) the whole policy network can be trained end-to-end with signals directly coming from actual actions.

In Figure 1, we illustrate the whole pipeline of our master-slave multi-agent architecture. Specifically, we demonstrate the network structure unfolded at two consecutive time steps. In the left part, for example, at time step $t$, the state $\boldsymbol{s}$ consists of each $s^i$ of the $i$-th slave agent and $s^m = \boldsymbol{o}_t$ of the master agent. Each slave agent is represented as a blue circle and the master agent is represented as a yellow rectangle. All agents are policy networks realized with RNN modules such as LSTM Hochreiter & Schmidhuber (1997) cells or a stack of RNN/LSTM cells. Therefore, besides the states, all agents also take the hidden state of RNN $\boldsymbol{h}_{t-1}$ as their inputs, representing their reasoning along time. Meanwhile the master agent also take as input some information from each slave agent $\boldsymbol{c}_i$ and broadcasts back its action output to all agents to help forming their final actions. These communications are represented via colored connections in the figure.

To merge the actions from the master agent and those from the slave agents, we propose a gated composition module (GCM), whose behavior resembles LSTM. Figure 4 illustrates more details. Specifically, this module takes the "thoughts" or hidden states of the master agent $h_t^m$ and the slave agents $h_t^i$ as input and outputs action proposals $a_t^{m \to i}$, which later will be added to independent action proposals from the corresponding slave agents $a_t^i$. Since such a module depends on both the "thoughts" from the master agent and that from certain slave agent, it facilitates the master to provide different action proposals to individual slave agents. We denote this solution as "MS-MARL+GCM". In certain cases, one may also want the master to provide unified action proposals to all agents. This could easily be implemented as a special case where the gate related to the slave's "thoughts" shuts down, which is denoted as regular MS-MARL.

## 3.2 LEARNING STRATEGY

As mentioned above, due to our design, learning can be performed in an end-to-end manner by directly applying policy gradient in the centralized perspective. Specifically, one would update all parameters following the policy gradient theorem Sutton & Barto (1998) as

$$\theta \leftarrow \theta + \lambda \sum_{t=1}^{T-1} \nabla_\theta \log \pi_\theta(\boldsymbol{s}_t, \boldsymbol{a}_t) v_t \tag{1}$$

where data samples are collected stochastically from each episode $\{\boldsymbol{s}_1, \boldsymbol{a}_1, r_2, ..., \boldsymbol{s}_{T-1}, \boldsymbol{a}_{T-1}, r_T\} \sim \pi_\theta$ and $v_t = \sum_{j=1}^t r_t$. Note that, for discrete action space, we applied softmax policy on the top layer and for continous action space, we adopted Gaussian policy instead. We summarize the whole algorithm in Algorithm 1 as follows.

---

**Algorithm 1:** Master-Slave Multi-Agent Policy Gradient (One Batch)

---

Randomly initialize $\theta = (\theta^m, \theta^1, \ldots, \theta^C)$;
Set $\sigma = 0.05$, $R_i = 0$;
Set $s_1 =$initial state, $o_1 =$initial occupancy map, $t = 0$;
**for** *b=1,BatchSize* **do**
    **while** $s_t \neq$*terminal* **and** $t < T$ **do**
        $t = t + 1$;
        **for** *each slave agent* $i = 1 \ldots C$ **do**
            Observe local state of agent $i$: $s_t^i$;
            Feed-forward: $(h_t^i, h_t^m) = MSNet(o_t, s_t^i, h_{t-1}^i, h_{t-1}^m; \theta^i, \theta^m)$;
            Sample action $a_i^t$ according to softmax policy or Gaussian policy;
            Execute action $a_i^t$;
            Observe reward $r_t^i$ (e.g. according to (2));
            Accumulate rewards $R_i = R_i + r_t^i$;
        Observe the next state $s_{t+1}$ and occupancy map $o_{t+1}$;
    Compute terminal reward $r_{terminal}$ (e.g. using (3));
    $r_{terminal} = \begin{cases} 1 & if \; battle \; won \\ -0.2 & else \end{cases}$ ;
    $R_i = R_i + r_{terminal}, \forall i = 1 \ldots C$;
Update network parameter $\theta^m, \theta^1, \ldots, \theta^C$ using (1);

---

## 3.3 DISCUSSIONS

As stated above, our MS-MARL proposal can leverage advantages from both the centralized perspective and the decentralized perspective. Comparing with the latter, we would like to argue that, not only does our design facilitate regular communication channels between slave agents as in previous works, we also explicitly formulate an independent master agent reasoning based on all slave agents' messages and its own state. Later we empirically verify that, even when the overall information revealed does not increase per se, an independent master agent tend to absorb the same information within a big picture and effectively helps to make decisions in a global manner. Therefore
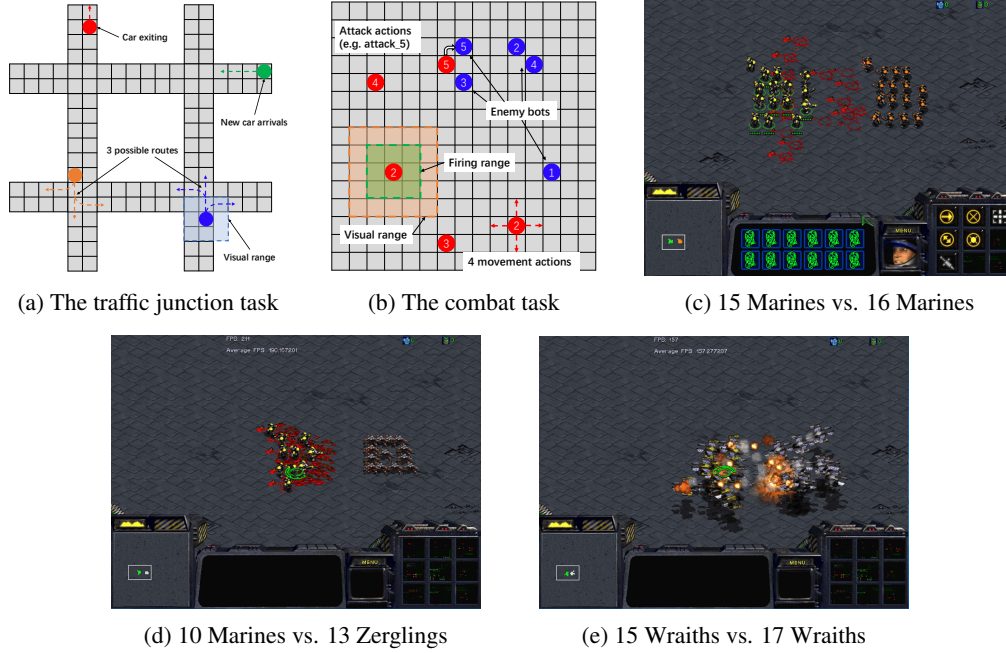
(a) The traffic junction task        (b) The combat task        (c) 15 Marines vs. 16 Marines

(d) 10 Marines vs. 13 Zerglings        (e) 15 Wraiths vs. 17 Wraiths

Figure 2: Environments for evaluation: (a) and (b) are established upon MazeBase EnvironmentSukhbaatar et al. (2015) and originally proposed in Sukhbaatar et al. (2016), (c)(d)(e) are minigames from StarCraft that is included in the TorchCraft platformSynnaeve et al. (2016)

compared with pure in-between-agent communications, MS-MARL is more efficient in reasoning and planning once trained.

On the other hand, when compared with methods taking a regular centralized perspective, we realize that our master-slave architecture explicitly explores the large action space in a hierarchical way. This is so in the sense that if the action space is very large, the master agent can potentially start searching at a coarse scale and leave the slave agents focus their efforts in a more fine-grained domain. This not only makes training more efficient but also more stable in a similar spirit as the dueling Q-network design Wang et al. (2015), where the master agent works as base estimation while leaving the slave agents focus on estimating the advantages. And of course, in the perspective of applying hierarchy, we can extend master-slave to master-master-slave architectures etc.

## 4    EXPERIMENTS

### 4.1    EVALUATION ENVIRONMENTS

To justify the effectiveness of the proposed master-slave architecture, we conducted experiments on five representative multi-agent tasks and environments in which multiple agents interact with each other to achieve certain goals. These tasks and environments have been widely used for the evaluation of popular MARL methods such as Sukhbaatar et al. (2016); Mao et al. (2017); Peng et al. (2017); Usunier et al. (2016); Foerster et al. (2017a;b).

The first two tasks are the traffic junction task and the combat task both originally proposed in Sukhbaatar et al. (2016). (Examples are shown in Figure 2(a)(b)) These two tasks are based on the MazeBase environment Sukhbaatar et al. (2015) and are discrete in state and action spaces. Detailed descriptions are given below:

**The traffic junction task**    As originally designed in Sukhbaatar et al. (2016), there are three versions of this task corresponding to different hard levels. In our experiments, we choose the hardest version which consists of four connected junctions of two-way roads (as shown in Fig. 2 (a)). During every time step, new cars enter the grid area with a certain probability from each of the eight

directions. Each car occupies a single-cell in the grid. All the cars keep to the right-hand side of the road. Note that there are only three possible routes for each car. A car is allowed to perform two kinds of actions at each time step: advances by one cell while keeping on its route or stay still at the current position. Once a car moves outside of the junction area, it will be removed. A collision happens when two cars move to the same location. In our experiments, following the literature, we apply the same state and reward originally designed in Sukhbaatar et al. (2016).

**The combat task** This environment simulates a simple battle between two opposing teams. The entire area is a $15 \times 15$ grid as shown in Fig. 2 (b). Both teams have 5 members and will be born at a random position. Possible actions for each agent are as follow: move one step towards four directions; attack an enemy agent within its firing range; or keep idle. One attack will cause an agent lose 1 point of health. And the initial health point of each agent is set to 3. If the health point reduces to 0, the agent will die. One team will lose the battle if all its members are dead, or win otherwise. Default settings from Sukhbaatar et al. (2016) are used in our experiments following the literature.

The following StarCraft micromanagement tasks originate from the well-known real-time strategy (RTS) game: StarCraft and was originally defined in Synnaeve et al. (2016); Usunier et al. (2016). Instead of playing a full StarCraft game, the micromanagement tasks involve a local battle between two groups of units. Similar to the combat task, the possible actions also include move and attack. One group wins the game when all the units of the other group are eliminated. However, one big difference from the combat task is that the StarCraft environment is continuous in state and action space, which means a much larger search space for learning. In other words, the combat task can also be considered as a simplified version of StarCraft micromanagement task with discrete states and actions.
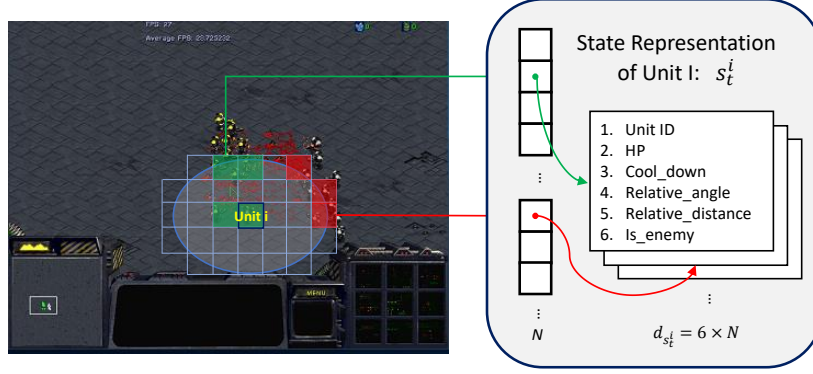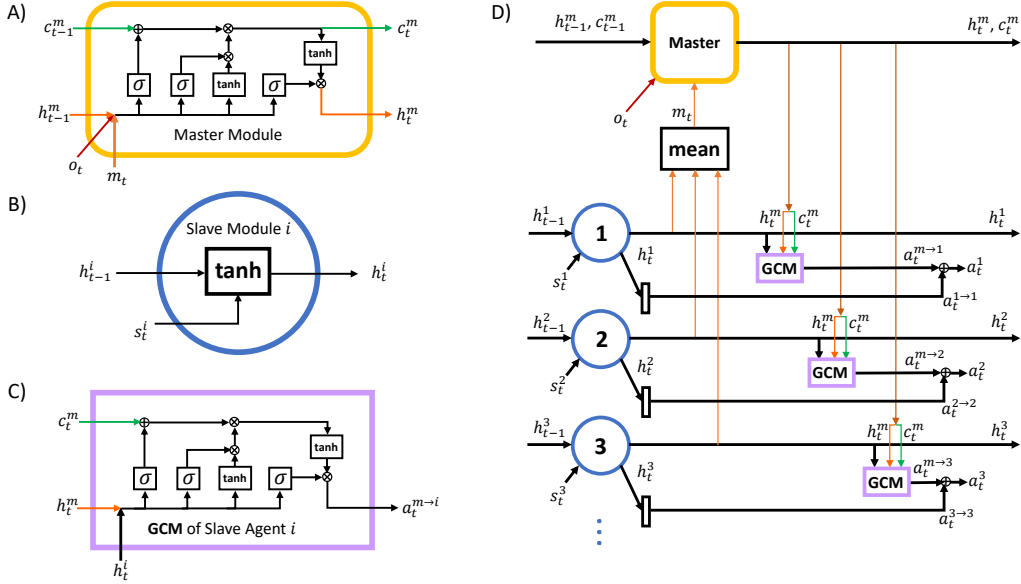
The selected micromanagement tasks for our experiments are {15 Marines vs. 16 Marines, 10 Marines vs. 13 Zerglings, 15 Wraiths vs. 17 Wraiths} (shown in Figure 2 (c)-(e)). All the three tasks are categorized as "hard combats" in Peng et al. (2017). Thus it is quite hard for an AI bot to win the combat without learning a smart policy. On all the three tasks, independent policy methods have been proved to be less effective. Instead, MARL methods must learn effective collaboration among its agents to win these tasks e.g. Peng et al. (2017); Usunier et al. (2016).

**15 Marines vs. 16 Marines** In this task, one needs to control 15 Terran marines to defeat a built-in AI that controls 16 Terran marines (showcased in Figure 2(a)). Note that the Marines are ranged ground units. Among the "mXvY" tasks defined in Synnaeve et al. (2016), such as m5v5 (5 Marines vs. 5 Marines), m10v10 (10 Marines vs. 10 Marines) m18v18 (18 Marines vs. 18 Marines), the chosen combat "15 Marines vs. 16 Marines" represents a more challenging version in the sense that the total number of agents is high and that the controlled team has less units than the enemies. Therefore a model need to learn very good strategies to win the battle. As described in Usunier et al. (2016), the key to winning this combat is to focus fire while avoiding "overkill". Besides this, another crucial policy - "Spread Out" is also captured in our experiments.

**10 Marines vs. 13 Zerglings** While we still control Marines in this scenario, the enemies belong to another type of ground unit - Zerglings. Unlike Terran Marines, Zerglings can only attack by direct contact, despite their much higher moving speed. Useful strategies also include "focus fire without overkill", which is similar to the "15 Marines vs. 16 Marines" task. Interestingly, unlike the "15M vs. 16M" task, the "Spread Out" strategy is not effective anymore in this case according to our experiments.

**15 Wraiths vs. 17 Wraiths** As a contrast to the above two tasks, this one is about Flying units. 15 Wraiths (ranged flying unit) are controlled to fight against 17 opponents of the same type. An important difference from the ground units is that the flying units will not collide with each other. Hence it is possible for flying units to occupy the same tile. In this case, the "Spread Out" strategy may not be important anymore. But it is essential to avoid "overkill" in this task since Wraiths have much longer "cooldown" time and much higher attack damage.

As presented in Peng et al. (2017), the state-of-the-art performance on these three tasks are still quite low compared with others. Amazingly, the proposed MS-MARL method achieves much higher winning rates, as demonstrated in later section. We will refer to these three tasks as {15M vs. 16M, 10M vs. 13Z, 15W vs. 17W} for brevity in the rest of this paper.

Figure 3: State definition of unit i in the task of {15M vs. 16M, 10M vs. 13Z, 15W vs. 17W}



Figure 4: A) master module, B) slave module, C) gated composition module of slave i, and D) specific model architecture.

## 4.2 IMPLEMENTATION DETAILS

In order to help reproduce our results, we hereby provide our very implementation details.

**Model Architecture** A detailed illustration of our implementation is displayed in Figure 4. As introduced in section 3, in our MS-MARL model, both the master agents and the slave agents are implemented using RNN or LSTM to incorporate temporal recurrence. Here in Figure 4, we use LSTM for the master module, and RNN for the slave module. The dimension of the hidden states in RNN or LSTM (including cell states) are all set to 50 for both the master and slave agents. The GCM module in Figure 4 (c) is noted as the "Gated Composition Module", which is introduced in section 3.1 in detail. Note that the action output is different for discrete and continuous tasks. For the traffic junction task and the combat task, the output of the network is designed as the probability of a number of actions since the action space is discrete. As a contrast, for "15M vs. 16M" "10M vs. 13Z" and "15W vs. 17W", our network directly generates a continuous action following Gaussian policy as described in section 3.

**State Features**   For the traffic junction task and the combat task, we just apply the original state features designed in Sukhbaatar et al. (2016). For "15M vs. 16M", "10M vs. 13Z" and "15W vs. 17W", we adopt one representation similar to existing methods Peng et al. (2017); Foerster et al. (2017a;b). To be specific, the feature of an agent is a map of a certain-sized circle around it, as illustrated in Figure 3. If there is an agent in a cell, a feature vector shown in Fig. 3 will be included in the final observation map.

**Action Definition**   For action definitions of the first two tasks, we use the default ones from Sukhbaatar et al. (2016). Since "15M vs. 16M", "10M vs. 13Z" and "15W vs. 17M" needs a continuous action space. We apply the design from Peng et al. (2017). It contains a 3-dimension vector, each of which is of the range from -1 to 1, representing action type, relative angle and distance to the target respectively.

**Reward Definition**   The reward definition of the traffic junction task and the combat task is the same as in Sukhbaatar et al. (2016). Here we define the reward for "15M vs. 16M", "10M vs. 13Z" and "15W vs. 17W" by simulating the design of that in the combat task. For each time step, there is a reward formulated in 2 to evaluate how the current action of agent i works.

$$r_t^i = \Delta n_{j \in N_m(i)}^t - \Delta n_{k \in N_e(i)}^t \tag{2}$$

Note that here $\Delta n_{j \in N_m(i)}^t = n_{j \in N_m(i)}^t - n_{j \in N_m(i)}^{t-1}$ (the definition of $\Delta n_{k \in N_e(i)}^t$ is similar) which is actually the change of the number of our units in the neighboring area. $N_m(i)$ and $N_e(i)$ refers to the sets containing our units and enemy units within a certain range of the current agent i. In experiments we define the range as a circle of a radius of 7 around agent i, which is exactly the circular range as 3 shows. Besides, we also define a terminal reward (as formulated in (3)) to add to the final reward at the end of each episode.

$$r_{terminal} = \begin{cases} 1 & if\ battle\ won \\ -0.2 & else \end{cases} \tag{3}$$

**Training Procedure** The entire training process is shown in Algorithm 1. Note that we follow the softmax policy as proposed in Sukhbaatar et al. (2016) for the first two discrete tasks. As for the cotinuous tasks {15M vs. 16M, 10M vs. 13Z, 15W vs. 17W}, a Gaussian policy is adopted for stochastic policy parameterization. Specifically, the output of the network is considered as the mean of a Gaussian policy distribution. By fixing the variance to $\sigma = 0.05$, the final actions are sampled as follows

$$a_i^t = \mu_t^i(s_t^i; \theta) + \sigma \cdot \mathcal{N}(0, 1), \tag{4}$$

And the score function in (1) can be computed as

$$\nabla_\theta \log \pi_\theta(s_t^i, a_t^i) = \frac{(a_t^i - \mu(s_t^i; \theta))}{\sigma^2} \frac{\partial \mu(s_t^i; \theta)}{\partial \theta}. \tag{5}$$

For the traffic junction task, we use a batch size of 16. For the combat task, a larger batch size of 144 is adopted for the training of both CommNet and our master-slave models. And for {15M vs. 16M, 10M vs. 13Z, 15W vs. 17W}, we find that a small batch size of 4 is good enough to guarantee a successful training. The learning rate is set to 0.001 for the first two tasks and 0.0005 for {15M vs. 16M, 10M vs. 13Z, 15W vs. 17W}. The action variance for {15M vs. 16M, 10M vs. 13Z, 15W vs. 17W} is initialized as 0.01 and drops down gradually as the training proceeds. For all tasks, the number of batch per training epoch is set to 100.

**Baselines**   For comparison, we select three state-of-the-art MARL methods that have been tested on these tasks of interest. A brief introduction of these baselines are given as follow:

- **CommNet:**   This method exploits a simple strategy for multi-agent communication. The idea is to gather a message output from all agents and compute the mean value. Then the mean value is spread to all agents again as an input signal of the next time step. In this way, a communication channel is established for all agents. Such a simple strategy has been shown to work reasonably well for the first two tasks in Sukhbaatar et al. (2016) and for {15M vs. 16M, 10M vs. 13Z, 15W vs. 17W} when explored in Peng et al. (2017).

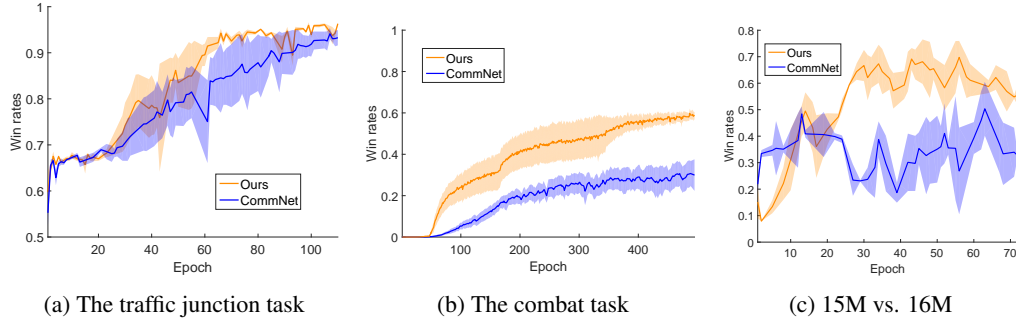(a) The traffic junction task    (b) The combat task    (c) 15M vs. 16M

Figure 5: Comparing winning rates of different methods on all three tasks

- **GMEZO:**   The full name of this method is GreedyMDP with Episodic Zero-Order Optimization. And it was proposed particularly to solve StarCraft micromanagement tasks in Usunier et al. (2016). The main idea is to employ a greedy update over MDP for inter-agent communications, while combining direct exploration in the policy space and backpropagation.

- **BiCNet:**   A bidrectional RNN is proposed to enable communication among multiple agents in Peng et al. (2017). The method adopts an Actor-Critic structure for reinforcement learning. It has been shown that their model is able to handle several different micromanagement tasks successfully as well as learn meaningful polices resembling skills of expert human players.

## 4.3   PERFORMANCE

Table 1 and Table 2 demonstrate the performance improvement of our method when compared with the baselines. For CommNet we directly run the released code on the traffic junction task and the combat task using hyper-parameters provided in Sukhbaatar et al. (2016). We compute the mean winning rates in Table 1 by testing the trained models for 100 rounds. However, since the code of GMEZO and BiCNet is not released yet, there is no report of their performance on traffic junction and combat tasks. Therefore we only compare with CommNet on these two tasks. And it can be seen that our MS-MARL model performs better than CommNet on both of the two tasks.

On the more challenging StarCraft micromanagement tasks {15M vs. 16M, 10M vs. 13Z, 15W vs. 17W}, the mean winning rates of GMEZO, CommNet and BiCNet are all available from Peng et al. (2017). For the performance of GMEZO, we directly follow results reported in Peng et al. (2017) since we utilize a similar state and action definition as them, but different from the original GMEZO Usunier et al. (2016). The results are displayed in Table 2. Obviously, on all three tasks, our MS-MARL method achieves consistently better performance compared to the available results of the three baselines.

A particular interesting phenomenon is about the effect of GCM. As can be seen from Table 2, only in the task of "15M vs. 16M" does GCM lead to a 5% improvement compared with the regular version of MS-MARL (82% vs. 77%). Strangely, there seems to be no benefit from GCM on the rest two micromanagement tasks - "10M vs. 13Z" and "15W vs. 17W". Intuitively, the main benefit of GCM is to allow different messages from the master to the slaves. Thus we can see it as enabling the master to encourage slaves to take the "spread out" action. In the case of "15M vs. 16M", such "spread out" strategy is critical since it allows our Marines to quickly bypass units standing in its way. In this case, it is much easier for our marines to attack the frontier enemies and eliminate them before the arrival of others (as discussed in section 4.5). As a contrast, the particular characteristics of 10M vs. 13Z and 15W vs. 17W make such strategy useless. In the specific case of "10M vs. 13Z", the Zerglings will always rush to the Marines rapidly. Thus they are always within the attack range of Marines during the whole battle. That is to say, the benefit of "Spread Out" - an early heavy attack on the frontier enemies, no longer exists. In such case, it is no wonder why "Spread Out" does not work anymore. The same situation is for the "15W vs. 17W" task, where the flying units never

10

collide with each other and can occupy the same tile. Thus all the units can always be attacked by the opponents, just as in the case of "10M vs. 13Z".

| Tasks | CommNet | MS-MARL |
|---|---|---|
| Traffic Junction | 0.94 | **0.96** |
| Combat | 0.31 | **0.59** |

Table 1: Mean winning rates of different methods on the first two discrete tasks

| Tasks | GMEZO | CommNet | BiCNet | MS-MARL | MS-MARL + GCM |
|---|---|---|---|---|---|
| 15M vs. 16M | 0.63 | 0.68 | 0.71 | 0.77 | **0.82** |
| 10M vs. 13Z | 0.57 | 0.44 | 0.64 | 0.75 | **0.76** |
| 15W vs. 17W | 0.42 | 0.47 | 0.53 | **0.61** | 0.60 |

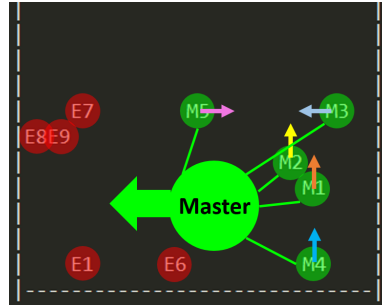Table 2: Mean winning rates of different methods on the cotinuous StarCraft micromanagement tasks

Figure 5 shows the training process of CommNet and our MS-MARL method by plotting winning rate curves for the first two tasks as well as "15M vs. 16M". As analyzed above, one key difference between our MS-MARL model and CommNet is that we explicitly facilitate an independent master agent reasoning with its own state and messages from all slave agents. From this plot, our MS-MARL model clearly enjoys better, faster and more stable convergence, which highlights the importance of such a design facilitating independent thinking.

## 4.4 ABLATION ANALYSIS

In the setting of the combat task, we further analyze how different components of our proposal contribute individually. Specifically, we compare the performance among the CommNet model, our MS-MARL model without explicit master state (e.g. the occupancy map of controlled agents in this case), and our full model with an explicit occupancy map as a state to the master agent. As shown in Figure 7 (a)(b), by only allowing an independently thinking master agent and communication among agents, our model already outperforms the plain CommNet model which only supports broadcasting communication of the sum of the signals. Further more, by providing the master agent with its unique state, our full model finally achieves a significant improvement over the CommNet model. Note here that, although we explicitly input an occupancy map to the master agent, the actual information of the whole system remains the same. This is because every information revealed from the extra occupancy map is by definition included to each agents state as their positions.



(a) Ablation of winning rates      (b) Master and slave actions

Figure 6: Ablation with occupancy map as master's input

Another interesting observation is how the master agent and each of the slave agents contribute to the final action choices (as shown in Fig. 6). We observe that the master agent does often learn an effective global policy. For example, the action components extracted from the master agent lead the

whole team of agents move towards the enemies regions. Meanwhile, all the slave agents adjust their positions locally to gather together. This interesting phenomenon showcases different functionality of different agent roles. Together, they will gather all agents and move to attack the enemies, which seems to be quite a reasonable strategy in this setting.

## 4.5 ANALYSIS OF LEARNED POLICIES

Note that the CommNet model has already been verified to have learned meaningful policies in these tasks according to Sukhbaatar et al. (2016); Peng et al. (2017). However, in our experiments, we often observe more successful policies learned by our method which may not be captured very well by the CommNet model.

For example, in the case of combat task, we often observe that some agents of the CommNet model just fail to find the enemies (potentially due to their limited local views and ineffective communications) and therefore lose the battle in a shorthanded manner, see e.g. Figure 7 (a). As a contrast, the learned model of our MS-MARL method usually gathers all agents together first and then moves them to attack the enemies. Figure 7 (b) showcases one example.

Another support case is in the task of "15M vs. 16M". In this task, our model learns a particular policy of spreading the agents into a half-moon shape ("Spread Out") to focus fire and attacking the frontier enemies before the others enter the firing range, as demonstrated in Fig. 8 (b). Actually, this group behavior is similar to the famous military maneuver "Pincer Movement" which is widely exploited in representative battles in the history. Although CommNet sometimes also follows such kind of policy, it often fails to spread to a larger "pincer" to cover the enemies and therefore loses the battle. Figure 8 (a) shows one of such examples. The size of the "pincer" seems especially important for winning the task of "15M vs. 16M" where we have less units than the enemies.
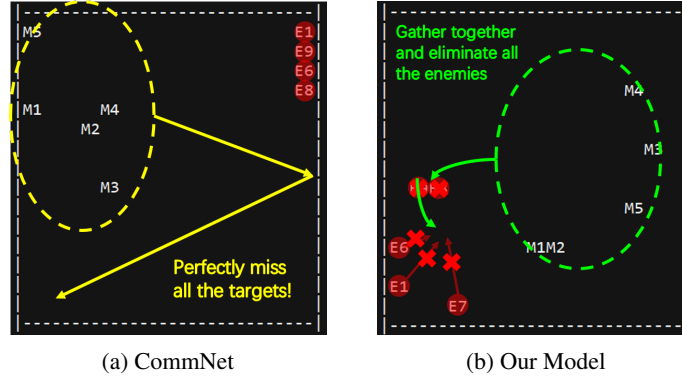


| (a) CommNet | (b) Our Model |

Figure 7: Two comparative cases in Combat Task: (a) a failure case that CommNet misses the targets (b) a successful case of our MS-MARL model
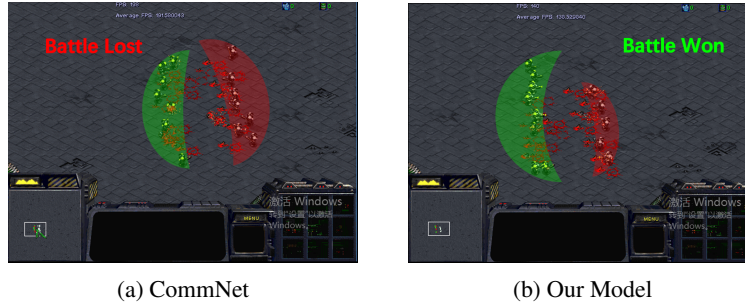


| (a) CommNet | (b) Our Model |

Figure 8: Illustration of the policy learned by our MS-MARL method: (a) is a failure case of CommNet (b) showcases the successful "Pincer Movement" policy learned by our MS-MARL model

## 5 CONCLUSION

In this paper, we revisit the master-slave architecture for deep MARL where we make an initial stab to explicitly combine a centralized master agent with distributed slave agents to leverage their individual contributions. With the proposed designs, the master agent effectively learns to give high-level instructions while the local agents try to achieve fine-grained optimality. We empirically demonstrate the superiority of our proposal against existing MARL methods in several challenging mutli-agent tasks. Moreover, the idea of master-slave architecture should not be limited to any specific RL algorithms, although we instantiate this idea with a policy gradient method, more existing RL algorithms can also benefit from applying similar schemes.

## REFERENCES

Lucian Busoniu, Robert Babuska, and Bart De Schutter. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, And Cybernetics-Part C: Applications and Reviews, 38 (2), 2008*, 2008.

Yongcan Cao, Wenwu Yu, Wei Ren, and Guanrong Chen. An overview of recent progress in the study of distributed multi-agent coordination. *IEEE Transactions on Industrial informatics*, 9(1): 427–438, 2013.

Jakob Foerster, Yannis M Assael, Nando de Freitas, and Shimon Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 2137–2145, 2016.

Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. *arXiv preprint arXiv:1705.08926*, 2017a.

Jakob Foerster, Nantas Nardelli, Gregory Farquhar, Philip Torr, Pushmeet Kohli, Shimon Whiteson, et al. Stabilising experience replay for deep multi-agent reinforcement learning. *arXiv preprint arXiv:1702.08887*, 2017b.

Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.

Xiangyu Kong, Bo Xin, Yizhou Wang, and Gang Hua. Collaborative deep reinforcement learning for joint object search. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Advances in Neural Information Processing Systems*, pp. 3675–3683, 2016.

Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(39):1–40, 2016.

Yuxi Li. Deep reinforcement learning: An overview. *arXiv preprint arXiv:1701.07274*, 2017.

Hangyu Mao, Yan Ni, Zhibo Gong, Weichen Ke, Chao Ma, Yang Xiao, Yuan Wang, Jiakang Wang, Quanbin Wang, Xiangyu Liu, et al. Accnet: Actor-coordinator-critic net for" learning-to-communicate" with deep multi-agent reinforcement learning. *arXiv preprint arXiv:1706.03235*, 2017.

Laëtitia Matignon, Laurent Jeanpierre, Abdel-Illah Mouaddib, et al. Coordinated multi-robot exploration under communication constraints using decentralized markov decision processes. In *AAAI*, 2012.

Dalila B Megherbi and Manuel Madera. A hybrid p2p and master-slave architecture for intelligent multi-agent reinforcement learning in a distributed computing environment: A case study. In *Computational Intelligence for Measurement Systems and Applications (CIMSA), 2010 IEEE International Conference on*, pp. 107–112. IEEE, 2010.

DB Megherbi and Minsuk Kim. A hybrid p2p and master-slave cooperative distributed multi-agent reinforcement learning technique with asynchronously triggered exploratory trials and clutter-index-based selected sub-goals. In *Computational Intelligence and Virtual Environments for Measurement Systems and Applications (CIVEMSA), 2016 IEEE International Conference on*, pp. 1–6. IEEE, 2016.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

Kui-Hong Park, Yong-Jae Kim, and Jong-Hwan Kim. Modular q-learning based multi-agent cooperation for robot soccer. *Robotics and Autonomous Systems*, 35(2):109–122, 2001.

Peng Peng, Quan Yuan, Ying Wen, Yaodong Yang, Zhenkun Tang, Haitao Long, and Jun Wang. Multiagent bidirectionally-coordinated nets for learning to play starcraft combat games. *arXiv preprint arXiv:1703.10069*, 2017.

Zhou Ren, Xiaoyu Wang, Ning Zhang, Xutao Lv, and Li-Jia Li. Deep reinforcement learning-based image captioning with embedding reward. *arXiv preprint arXiv:1704.03899*, 2017.

David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.

Sainbayar Sukhbaatar, Arthur Szlam, Gabriel Synnaeve, Soumith Chintala, and Rob Fergus. Mazebase: A sandbox for learning from games. *arXiv preprint arXiv:1511.07401*, 2015.

Sainbayar Sukhbaatar, Rob Fergus, et al. Learning multiagent communication with backpropagation. In *Advances in Neural Information Processing Systems*, pp. 2244–2252, 2016.

Richard S Sutton and Andrew G Barto. *Introduction to reinforcement learning*, volume 135. MIT Press Cambridge, 1998.

Gabriel Synnaeve, Nantas Nardelli, Alex Auvolat, Soumith Chintala, Timothée Lacroix, Zeming Lin, Florian Richoux, and Nicolas Usunier. Torchcraft: a library for machine learning research on real-time strategy games. *arXiv preprint arXiv:1611.00625*, 2016.

Nicolas Usunier, Gabriel Synnaeve, Zeming Lin, and Soumith Chintala. Episodic exploration for deep deterministic policies: An application to starcraft micromanagement tasks. *arXiv preprint arXiv:1609.02993*, 2016.

Katja Verbeeck, Ann Nowé, and Karl Tuyls. Coordinated exploration in multi-agent reinforcement learning: an application to load-balancing. In *Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems*, pp. 1105–1106. ACM, 2005.

Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. FeUdal networks for hierarchical reinforcement learning. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 3540–3549, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. URL `http://proceedings.mlr.press/v70/vezhnevets17a.html`.

Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Van Hasselt, Marc Lanctot, and Nando De Freitas. Dueling network architectures for deep reinforcement learning. *arXiv preprint arXiv:1511.06581*, 2015.

Erfu Yang and Dongbing Gu. Multiagent reinforcement learning for multi-robot systems: A survey. Technical report, tech. rep, 2004.

Dayong Ye, Minjie Zhang, and Yun Yang. A multi-agent framework for packet routing in wireless sensor networks. *sensors*, 15(5):10026–10047, 2015.