

Learning Selective Communication for Multi-Agent Path Finding

Ziyuan Ma^{1*}, Yudong Luo^{2*} and Jia Pan³

Abstract—Learning communication via deep reinforcement learning (RL) or imitation learning (IL) has recently been shown to be an effective way to solve Multi-Agent Path Finding (MAPF). However, existing communication based MAPF solvers focus on broadcast communication, where an agent broadcasts its message to all other or predefined agents. It is not only impractical but also leads to redundant information that could even impair the multi-agent cooperation. A succinct communication scheme should learn *which* information is relevant and influential to each agent’s decision making process. To address this problem, we consider a request-reply scenario and propose *Decision Causal Communication* (DCC), a simple yet efficient model to enable agents to select neighbors to conduct communication during both training and execution. Specifically, a neighbor is determined as relevant and influential only when the presence of this neighbor causes the decision adjustment on the central agent. This judgment is learned only based on agent’s local observation and thus suitable for decentralized execution to handle large scale problems. Empirical evaluation in obstacle-rich environment indicates the high success rate with low communication overhead of our method.

I. INTRODUCTION

Multi-Agent Path Finding is the problem of arranging a set of collision-free paths for a set of agents on a given graph. Although MAPF is NP-hard to solve optimally on graphs [1] and even 2D grids [2], many optimal MAPF algorithms have been developed in recent years. Some reduce MAPF to other well-studied problems, e.g., ILP [3] and SAT [4], others solve it with search-based algorithms, e.g., enhanced A* search [5], [6], Conflict-Based Search [7] and its improved variants [8]. However, the limitation of these centralized planning methods is that they do not scale well to a large number of agents.

RL and IL methods with decentralized execution have been applied to address this issue [9], [10], [11], [12], [13], [14], [15]. During execution, each agent takes action based on its individual decision model with its own observation (may also include messages from other agents) as input, which avoids the scalability problem. For RL-based methods [9], [10], [12], [15], MAPF is generally modeled as a Markov Game [16] with partial observability, where a reward function is designed for each agent, and the goal of each agent is to maximize its expected total return. Usually, expert guidance is applied to guide the RL during training in terms of behavior cloning [9], shaped reward [12], or

heuristic map [15]. For IL-based methods [11], [13], [14], partial observability is still considered, but the objective is to minimize the divergence between expert demonstrations and generated trajectories, so that the learned policy can mimic the expert behavior.

Recently, researchers focus on trainable communication channels between agents for MAPF, where extra information can be obtained during both training and execution to enhance multi-agent cooperation [13], [14], [15]. However, these methods focus on broadcast communication in which messages will be transmitted to all other or predefined set of agents. For instance, in [13] and [14], the messages are broadcast to all other agents within a distance (communication radius). In [15], the central agent takes the messages from the two nearest neighbors after getting all the broadcast messages from its neighbors inside the field of view (FOV). Although empirical results have shown great improvements of communication based methods compared with non communication ones [15], the drawback is that broadcast communication requires lots of bandwidth and incurs additional system overhead and latency in practice. More importantly, not every agent provides useful information for cooperation, and redundant information can even impair the learning process.

Plenty of methods on reducing communication overhead or learning selective communication have been proposed in the literature of RL [17], [18], [19], [20], [21], [22]. But most of them are designed for centralized training framework on cooperative Markov Game, where agents share the same team reward. In contrast, each agent has its individual reward in MAPF, and thus more information needs to be swapped for cooperation. Therefore, selective communication is much more important in this case. To design a succinct communication protocol for MAPF, we consider a request-reply communication scenario where messages are updated in two stages, and propose *Decision Causal Communication*, a simple yet efficient mechanism to enable agents to select other members in the swarm for communication during both training and execution. More specifically, before communication, each agent is capable to figure out which neighbors in its FOV are relevant and influential by just using local observation. The agent makes two temporary decisions based on its raw local observation and the modified local observation by masking out a neighbor. For any neighbor that can cause difference between these two temporary decisions, that neighbor is considered as relevant and influential, and the agent will send a request to that neighbor for communication purpose. This communication protocol is also naturally capable for decentralized execution. Empirical results show the high suc-

¹School of Computing Science, Simon Fraser University, Canada ziyuan.ma@sfu.ca

²School of Computer Science, University of Waterloo, Canada yudong.luo@uwaterloo.ca

³Department of Computer Science, The University of Hong Kong, Hong Kong, China jpan@cs.hku.hk

* indicates equal contribution.

cess rate with low communication overhead of our method compared with its counterparts.

Contributions. In this work, we propose DCC for MAPF. Unlike broadcast communication, DCC learns selective communication, and hence encourages agents to focus only on relevant information. We demonstrate that, on one hand, selective communication can remove temporal redundant messages, which is beneficial to multi-agent cooperation. On the other hand, communication overhead is greatly cut down due to the reduction of communication frequency.

II. RELATED WORK

A. Reinforcement Learning based MAPF

RL-based planners generally formulate MAPF as a multi-agent reinforcement learning (MARL) problem to learn decentralized policies for agents (robots). Compared with centralized planning methods, learning decentralized policies scales well to a large number of agents. Collision free policies are usually learned by guiding RL with constraints or demonstration data [23], [9], [10], [12]. Incorporating demonstration guidance with RL can be divided into two categories, using demonstrations generated by a centralized planner, or by a single-agent planner. For instance, a well known framework named PRIMAL [9] builds on the asynchronous advantage actor critic (A3C) network [24] as its RL part and uses behaviour cloning to supervise the training of RL. So it requires demonstrations generated by a centralized planner named ODrM* [6]. The limitation of using a centralized planner is that it requires solving a MAPF problem and is thus time-consuming, especially in a complex environment with a large number of agents. For methods using single-agent planner as guidance, MAPPER [12] and Globally Guided RL (G2RL) [10] use A* search for single-agent path generation and apply an off-route penalty if agents fail to follow the path. The potential issue is that, as single-agent shortest paths are usually not unique and not globally optimal for multi-agent environment, forcing agents to follow these paths by extra shaped rewards can mislead agents.

To tackle the above issues, our previous method DHC [15] does not require a centralized planner. Although single-agent shortest paths are still adopted as guidance, DHC embeds all the potential choices of shortest paths as heuristic into the input of the model, instead of forcing agents to follow a specific path. Thus, no special shaped rewards are required.

Even though empirical results have shown that enhancing RL with guidance can help to learn collision free policies for MAPF, the cooperation among agents is not directly modeled. Communication via graph convolution is a promising way to achieve multi-agent cooperation, where an agent aggregates information from its neighbors, including itself, and makes decisions based on this augmented information [25], [26], [27]. This idea is recently deployed for MAPF. By treating each agent in the environment as a node and connecting neighboring nodes (agents) if they are inside the FOV of each other, a graph is formulated by DHC and multi-head dot-product attention [28] serves as the convolutional kernel to compute interactions among agents.

Similar ideas are applied by the work described in [13] and Message-Aware Graph Attention neTwork (MAGAT) [14], where the communication part is a graph neural network (GNN) [29]. However, it is worth noting that both [13] and [14] are IL-based methods, which train the model to imitate the demonstrations generated by Conflict-Based Search (CBS) [7] and Enhanced CBS [30], respectively. Thus, MAPF instances still need to be pre-solved while extreme scenarios may fail to be solved and collected into the training data. Also, none of these communication based methods consider the efficiency of communication, resulting in high communication overhead and latency.

B. Efficient and Selective Communication for MARL

Most existing works on communication in MARL, including those described above, focus on broadcast communication, i.e., broadcasting the messages to all other or predefined agents. To improve efficiency, a number of strategies have been proposed. Individualized Controlled Continuous Communication Model (IC3Net) [17] considers that full communication is not always necessary and agents can determine whether to send messages to others indicated by a gating function. However, this gating function can block all communication channels of an agent. Cases where an agent’s message is instructive to some agents but useless or harmful to others can not be handled. Schedule Communication (SchedNet) [18] only allows a limited number of agents to broadcast messages, who are chosen by some importance weights assigned. The shortcoming is that a central scheduler is required to gather all individual weights and decide which agents should be entitled to broadcast their messages. Variance Based Control (VBC) [19] lets agents transmit their messages only when the variance of the message vector is high (high variance implies that the message is informative). And agents decide whether to request messages based on local ambiguity. Temporal Message Control (TMC) [20] stores agent messages in the buffer and new messages are sent out only when they contain relatively new information compared with old ones. Nearly Decomposable Q-functions (NDQ) [21] optimizes communication via minimizing the entropy of messages between agents. Individually Inferred Communication (I2C) [22] learns a prior network, which takes local observation and identity of another agent as input, to predict a belief of whether to communicate. All VBC, TMC, NDQ and I2C are designed on the centralized training methods that factorize the joint action-value function, such as QMIX [31], different from MAPF settings.

III. PROBLEM FORMULATION

A. MAPF Problem Definition

MAPF has many variants as summarized in [32]. In this paper, we consider the classical MAPF case defined in [32] that (a) each agent performs an action in every time step, and may cause vertex and swapping conflicts, (b) agents “stay at target” until all agents have reached their goals, and (c) optimizes the sum of cost. This definition differs from the environments used in PRIMAL and MAPPER, which are

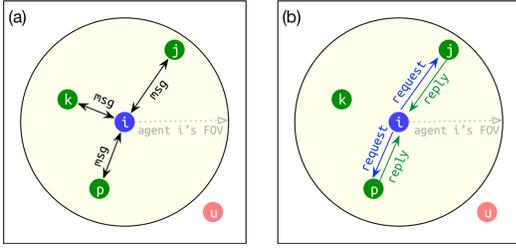


Fig. 1. (a) Broadcast communication: agent i broadcasts its message to all neighboring agents j , k , and p . And agents j , k , p will also broadcast their messages to agent i . If predefined scope is set, agent i can decide whose message to be used in calculation, e.g., DHC takes the nearest two agents. There is no communication between agent i and u , as agent u is outside agent i 's FOV. (b) Request-reply based selective communication: agent i omits irrelevant agent k , and only communicates with relevant agent j and p . After determining which agent to communicate, the communication is conducted in two rounds. At the first round, agent i sends out messages (requests) to agent j and p . At the second round, agent j and p reply with their updated messages to agent i .

much simpler. In PRIMAL, agents take actions in turn (only one agent takes action at a time), although we can treat this as all agents act simultaneously but only one agent can move while all others choose no-op action. In MAPPER, agents are removed from the environment upon reaching their goal locations, which simplifies the problem, as there will be no livelocks when an agent reaches its goal but obstructs other agents from getting to their goals. We also notice that the environment in [33] is specially designed where obstacles compose narrow corridors that only allow one agent to pass at a time. We do not consider an environment like this.

Formally, we define MAPF as follows. Given an undirected graph $G = (V, E)$ and a set of n agents, each agent indexed by i has a unique start vertex $s_i \in V$ and a unique goal vertex $g_i \in V$. Time is discretized into time steps. At each discrete time step, each agent can either *move* to an adjacent vertex or *wait* at its current vertex. A path for i -th agent is a sequence of adjacent (indicating a moving) or identical (indicating a waiting) vertices starting at s_i and terminating at g_i . Agents remain at their goal vertices (g_i) after they complete their paths. There are two kinds of collisions, either vertex collision or edge collision. A vertex collision is a tuple $\langle i, j, v, t \rangle$ where i -th and j -th agents reaching at the same vertex v at time t . An edge collision is a tuple $\langle i, j, u, v, t \rangle$ where i -th and j -th agents traverse the same edge (u, v) in opposite directions at time t . A solution to MAPF is a set of collision-free paths, one for each agent. The quality of a solution is measured by the sum of arrival time of all agents at their goal vertices.

B. Environment Setup

In line with standard MAPF tasks, we focus on 2D 4-neighbor grids where agents, goals, and obstacles occupy one grid cell respectively. Formally, the map is a $m \times m$ matrix, where 0 represents a free location and 1 is an obstacle. Each map is chosen n start positions and n corresponding goal positions for n agents from free locations. We make sure each goal is reachable from its start point and there is no overlap among $2n$ selected positions. At each time

TABLE I
REWARD FUNCTION DESIGN

Actions	Reward
Move (Up/Down/Left/Right)	-0.075
Stay (on goal, away goal)	0, -0.075
Collision (obstacle/agents)	-0.5
Finish	3

step, agents move simultaneously to neighboring locations or wait at their current locations. Thus, the size of the action space is 5 (move to four directions or stay still). Agents may hit obstacles or conflict with others during simulation. We handle the collision by recursively resetting the related agents to previous states until no collision exists.

We consider a partially observable environment, which is more compatible with the real world problem. Each agent can only access the information inside its FOV with size $\ell \times \ell$ ($\ell < m$), where ℓ is an odd number to make sure agents are at the center of FOV.

The reward design is adopted from DHC [15], shown in Table I. Generally, negative rewards are given to agents for each movement to facilitate goal reaching. No shaped reward is required in our method.

IV. SELECTIVE COMMUNICATION FOR MAPF

In this paper, we propose DCC¹, a simple yet efficient mechanism to reduce communication overhead in multi-agent systems. DCC can be instantiated by independent Q-learning [34], or any framework of centralized training and decentralized execution, such as QMIX [31] and VDN [35]. We consider a request-reply scenario different from the traditional setup where each agent only sends out an almost information-less indicator, e.g. a scalar value, as the request signal. Instead, in our setup, each agent sends out a request signal with rich information including its own messages along with the relative positions of neighbors. In this way, after receiving the request, the agents being requested can immediately benefit from this query by collecting some information from the query agent. Fig. 1 shows the difference between broadcast and our request-reply based selective communication. In the following part, We first introduce DCC communication mechanism and then instantiate DCC with independent Q-learning.

A. Decision Causal Communication

Our communication mechanism is motivated by I2C [22], where causal inference is utilized to derive a threshold to trigger communication. Specifically, I2C defines the causal effect of agent j to agent i as the KL-divergence [36] between two decision probabilities, $D_{KL}(P(a_i | \mathbf{a}_{-i}, \mathbf{o}) || P(a_i | \mathbf{a}_{-ij}, \mathbf{o}))$, where \mathbf{o} is the joint observations, \mathbf{a}_{-i} denotes the joint actions except for agent i , and similar meaning for \mathbf{a}_{-ij} . The magnitude of this divergence indicates how much agent i will adjust its policy if taken into account agent j 's policy. If

¹Code available at <https://github.com/ZiyuanMa/DCC>

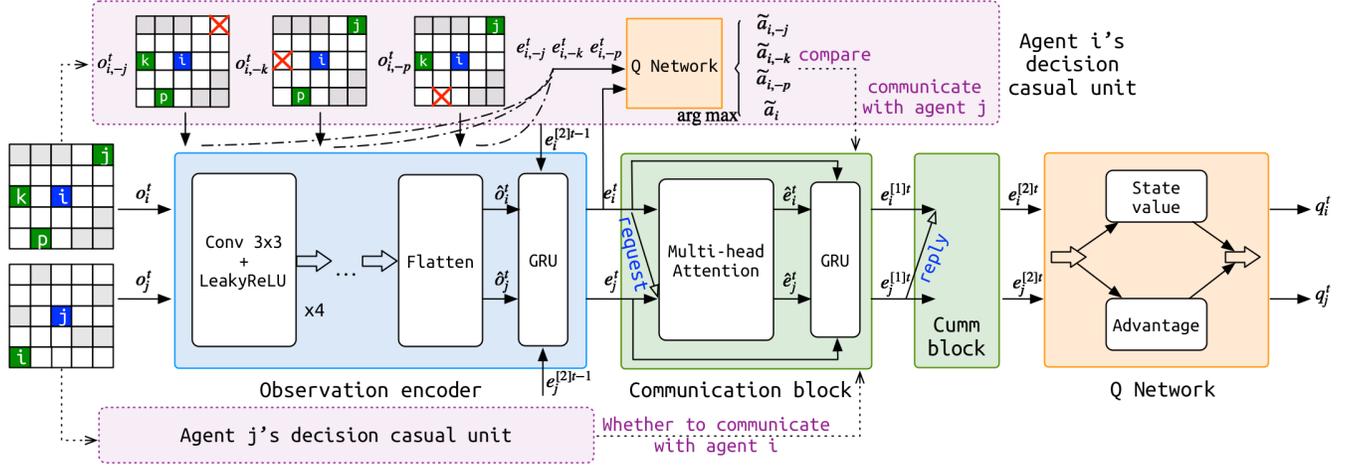


Fig. 2. System flow of DCC. It contains four modules: observation encoder (blue), decision casual unit (purple), communication block (green), and Q network (orange). An example of agent i 's FOV is shown by a 5×5 grid, where the blue color represents the agent itself, the green color represents neighbors, and the gray color represents obstacles. Agent i 's decision casual unit constructs three modified local observations by masking each neighbor inside the FOV, represented by a red cross. The embeddings by observation encoder for the raw and modified local observations are fed to Q network directly to get temporary decisions, which are further used to decide which neighbor is relevant for communication.

the divergence is small, agent i chooses not to communicate with j as that will hardly affect its own policy.

However, I2C is only effective for centralized training algorithms as it requires to know the joint observations and actions of all other agents. In real world problems such as MAPF, usually an agent with partial observability can only observe the existence of other agents instead of their policies. Although agents can broadcast their actions, that will definitely increase the communication overhead. In practice, if the policy of an agent requesting communication conditions on the action of the communicated agent, the circular dependencies can occur. At last, setting a reasonable threshold for that KL-divergence highly depends on empirical experiments and may vary in different problems.

Thus, to design an appropriate communication protocol for MAPF, in this work, we restrict the strategy in I2C by triggering communication between agent i and j only when the *existence* of agent j causes the policy adjustment on agent i , and we call this *Decision Causal Communication*. To get an intuition of the proposed protocol, consider the learning based MAPF methods discussed in Section II-A, where most of them guide the RL agent with global or local optimal demonstrations. In order to diminish the communication cost, an agent should follow the demonstration guided policy as often as possible without communicating with others, if the existence of other agents in its FOV will not affect its current policy. Otherwise, the agent communicates with those neighbors who will individually affect its policy. Although the policy may be adjusted after communication and we can not guarantee the adjusted policy is always better than unchanged policy during simulation, training the communication part along with RL in an end to end manner to maximize the expected total return will force the communication to be helpful and performed only when it is necessary.

Formally, this communication protocol can be formulated as follows. Suppose agent j can be seen by agent i (agent

j is inside the FOV of agent i), and the relative position of agent j in i 's FOV is (x_j, y_j) . Agent i gets its local observation o_i from the environment and constructs another modified observation $o_{i,-j}$ ($-j$ means without j) by setting the information at agent j 's position by some special value, such as zero: $o_{i,-j} \leftarrow o_i(x_j, y_j) = 0$. Then agent i uses its local decision model, without communication, to predict two temporary actions \tilde{a}_i and $\tilde{a}_{i,-j}$ based on o_i and $o_{i,-j}$, respectively. If these two actions match with each other, implying that the existence of agent j will not affect i 's policy, agent i will not request message from agent j . Otherwise, agent i sends a request to agent j via a communication channel, and agent j will respond with a message. In this work, we consider messages are updated in both request and reply stages, which is discussed in the following subsection. The final aggregated messages are further used by policy or Q-network to generate the final actions.

B. Deep Q-Learning Model Design

We instantiate DCC with independent Q-learning, namely the model is designed and trained from a single agent's perspective by treating other agents as part of the environment. We borrow some architecture design from DHC. The whole model of DCC consists of four main components, including observation encoder, decision causal unit, communication block, and Q-network. The model architecture is shown in Fig. 2. We begin from the model input and describe the four components one by one.

a) *Local observation input*: At time step t , agent i gets a matrix with the shape $\ell \times \ell \times 6$ as input, denoted by o_i^t . The first channel is a binary matrix representing the obstacles inside the FOV. The second channel is a binary matrix indicating the locations of other agents if within the FOV. The other four are heuristic channels proposed by DHC [15], to encode multiple choices of single-agent shortest path. We refer readers to the DHC paper for further details.

b) *Observation encoder*: The encoder contains four stacked convolutional layers with a GRU [37]. The local observation o_i^t is first encoded into \hat{o}_i^t by four convolutional layers. Then by adopting the last step communication outcome $e_i^{[2]t-1}$ as the hidden state, the GRU gets \hat{o}_i^t and generates the intermediate message e_i^t .

c) *Decision causal unit*: This unit depends on the observation encoder to embed modified observations, and the Q-network (will be discussed later), to get temporary actions. Let \mathbb{B}_i denote all neighbors of agent i . Given o_i^t , $o_{i,-j}^t$ is constructed according to Section IV-A for all $j \in \mathbb{B}_i$. Each $o_{i,-j}^t$ is passed through the observation encoder as discussed above to get an embedding $e_{i,-j}^t$. By skipping the communication, e_i^t and $\{e_{i,-j}^t\}_{j \in \mathbb{B}_i}$ are directly fed into the Q-network to compute Q-values. Actions \tilde{a}_i^t and $\{\tilde{a}_{i,-j}^t\}_{j \in \mathbb{B}_i}$ are inferred by applying argmax over Q-values. Based on these actions, the communication scope of agent i is

$$\mathbb{C}_i = \{j | \tilde{a}_i^t \neq \tilde{a}_{i,-j}^t\}_{j \in \mathbb{B}_i}. \quad (1)$$

Note that temporary actions \tilde{a}_i^t and $\{\tilde{a}_{i,-j}^t\}_{j \in \mathbb{B}_i}$ are only used to decide the communication scope, not the final action to be executed.

d) *Communication block*: This block is the graph convolution with multi-head dot-production as the convolutional kernel [28] followed by a GRU. We regard the request-reply as a two-round communication. After defining the communication scope \mathbb{C}_i , at the first round, agent i sends request information to all agent $j \in \mathbb{C}_i$. This information includes its own message e_i^t as well as the relative position of i 's neighbors inside i 's FOV, denoted as l_i^t . The relative position of each neighbor is originally represented as a one-hot vector, i.e. the vector length is $\ell \times \ell$, and further embedded by a neural network layer. At the first round, agent j may receive many requests from different agents, and we define this request receiving scope for agent j as

$$\bar{\mathbb{C}}_j = \{i | j \in \mathbb{C}_i\}. \quad (2)$$

Then each agent j who receives requests will update its message using multi-head attention. Let $\bar{\mathbb{C}}_{j+}$ denote the set $\{j, \bar{\mathbb{C}}_j\}$. For every agent $\bar{i} \in \bar{\mathbb{C}}_{j+}$, agent \bar{i} 's own message $e_{\bar{i}}^t$ is projected to **Query** by matrix \mathbf{W}_Q^h , while the concatenation of $e_{\bar{i}}^t$ and $l_{\bar{i}}^t$ is projected to **Key** and **Value** by matrix \mathbf{W}_K^h and \mathbf{W}_V^h in each independent attention head h . The relation between agent j and $\bar{i} \in \bar{\mathbb{C}}_{j+}$ in h -th attention head is computed as

$$\mu_{j\bar{i}}^h = \text{softmax} \left(\frac{\mathbf{W}_Q^h e_j^t \cdot (\mathbf{W}_K^h [e_{\bar{i}}^t, l_{\bar{i}}^t])^\top}{\sqrt{d_K}} \right), \quad (3)$$

where d_K is the dimension of **Keys**. The **Values** are weighted summed by weights $\mu_{j\bar{i}}^h$ over $\bar{\mathbb{C}}_{j+}$ at each head h . And all head outputs are concatenated over \mathcal{H} heads to pass through a neural network layer f_o for the final output

$$\hat{e}_j^t = f_o \left(\text{concat} \left(\sum_{\bar{i} \in \bar{\mathbb{C}}_{j+}} \mu_{j\bar{i}}^h \mathbf{W}_V^h [e_{\bar{i}}^t, l_{\bar{i}}^t], \forall h \in \mathcal{H} \right) \right). \quad (4)$$

The message \hat{e}_j^t and the initial message e_j^t are first aggregated by a GRU to generate the first round outcome $e_j^{[1]t}$. Then $e_j^{[1]t}$ acts as the initial message for the next round.

At the second round, agents who received requests will reply with their updated messages along with the relative position of their neighbors, i.e. agent i will receive $e_j^{[1]t}$ and l_j^t from all agent $j \in \mathbb{C}_i$. Let \mathbb{C}_{i+} denote the set $\{i, \mathbb{C}_i\}$. Agent i will update its message in the same manner as defined in Equations 3-4 by replacing j with i , and replacing $\bar{i} \in \bar{\mathbb{C}}_{j+}$ with $\bar{j} \in \mathbb{C}_{i+}$. Denote the message generated by Equations 3-4 at this round as $\hat{e}_i^{[1]t}$. Finally, $\hat{e}_i^{[1]t}$ and $e_j^{[1]t}$ are aggregated by a GRU to output $e_i^{[2]t}$, the final outcome message of the second round.

e) *Q-network*: The message $e_i^{[2]t}$ is adopted by a dueling Q-network [38], which separates state advantage $V_s(\cdot)$ and action advantage $A(\cdot)$, to predict the Q-values.

$$Q_{s,a}^i = V_s(e_i^{[2]t}) + \left(A(e_i^{[2]t})_a - \frac{1}{|\mathcal{A}|} \sum_{a'} A(e_i^{[2]t})_{a'} \right), \quad (5)$$

where \mathcal{A} is the action space.

After we get the Q values of agent i , the loss function for training the whole model is a time difference (TD) error. To stabilize training, we use a multi-step TD error calculated by the mean square error (MSE)

$$\mathcal{L}(\theta) = \text{MSE} \left(R_t^i - Q_{s_t, a_t}^i(\theta) \right) \quad (6)$$

with $R_t^i = r_t^i + \gamma r_{t+1}^i + \dots + \gamma^n Q_{s_{t+n}, a_{t+n}}^i(\bar{\theta})$, where r_t^i is the reward received by agent i at time t , R_t^i is the multi-step expected return, θ represents the parameters of the entail model, and $\bar{\theta}$ denotes the parameters of the target network, a periodical copy of the online parameters θ .

C. Training

As mentioned above, we train the model from a single agent's perspective using independent Q-learning, where each agent learns its own action value Q^i independently and simultaneously, treating other agents as part of the environment. The benefit is that it avoids the scalability problem in centralized training, which requires learning a Q-function for joint actions over all agents, as the joint action space grows exponentially when the number of agents increases. To further simplify the training process, instead of training separate models for different agents, we train a single model by sharing parameters among agents. To facilitate training, we adopt two training strategies, namely distributed training and curriculum learning [39], which are discussed in the following.

a) *Distributed training*: Distributed training has significantly improved the performance of RL compared with non distributed ones, such as R2D2 [40] and Ape-X [41]. In MAPF, PRIMAL is a distributed RL method by using A3C. The idea behind this distributed version of algorithm is to parallelize the gradient computation, to more rapidly optimize the parameters. Another approach, as proposed in Ape-X, is to parallelize experience data generation and selection with a shared replay memory. We adopt the latter

strategy, where we setup multiple runners on CPUs with their own copy of the environment and up to date model to generate experience data, and a single learner on GPU to train. The runners initialize the priorities for the experience data and feed into a global prioritized reply buffer. The learner samples the most useful experiences from the buffer to update the model parameters and also update the priorities of the experiences. The advantage is that the model training is allocated to GPU while original A3C is designed for multi-core CPU. Note that the transitions of all the agents need to be recorded for communication purpose in the model. As the priorities are shared among all runners, the good experiences explored by any runner can benefit the learner.

As we only have a single learner on GPU, the core task is to efficiently distribute multiple runners on CPUs. Here we utilize a powerful Python package named Ray [42], a distributed framework designed for machine learning, to easily deploy parallel programs on multiple CPUs with little modification to the deep learning code.

b) Curriculum learning: The ultimate goal of the learning model is to handle complex cases in a large environment with high obstacle density and many agents. However, it would be hard for the model to learn fast and stably if directly starting from a complicated training environment. As pointed in [39], presenting training examples not randomly but in a meaningful order will benefit machine learning algorithms. Motivated by this idea, we design a training pipeline to gradually introduce more complex training cases to the model. Specifically, training starts with a task where only 1 agent stays in a 10×10 environment. Then if the success rate of the current task exceeds 0.9, two new tasks are established for training. One is to increase the number of agents by 1, and the other is to increase the size of the environment by 5. The final task with 16 agents in a 40×40 environment will be reached as the training scale grows.

c) Other training settings: For the environment setting, the obstacle density for each environment during training is sampled from a triangular distribution between 0 and 0.5 with a peak at 0.33, which is the same as the obstacle setting in PRIMAL. The FOV size is 9×9 (10×10 in PRIMAL, we make it odd) and the agent is in the middle of the FOV. The maximum step for the environment is 256.

For training setting, we train the model with a batch size of 128 and a sequence length of 20. The discount factor $\gamma = 0.99$. The multi-step TD error defined in Equation 5 is computed with a length of 2. The number of runners for distributed training is set to be 16.

V. EXPERIMENTS

The model is trained and tested in classical MAPF environments as discussed in Section III. We compare DCC with one of the most well known RL baselines named PRIMAL [9], and the most recent RL plus communication method named DHC [15] in terms of success rate and average step. For comparing communication overhead, we build a baseline model, which uses request-reply mechanism but with predefined communication scope. PRIMAL and

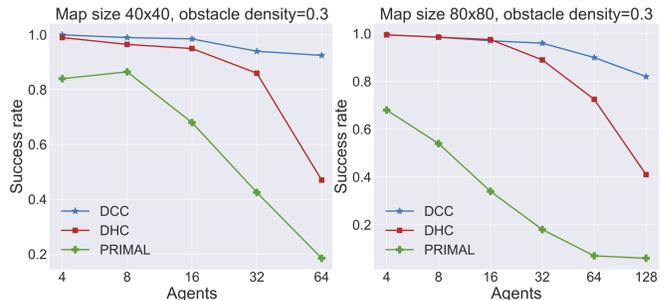


Fig. 3. Success rate of our method compared with PRIMAL and DHC in two different scenarios.

TABLE II

AVERAGE STEPS IN TWO TYPES OF ENVIRONMENTS WITH OBSTACLE DENSITY = 0.3

Average steps in 40×40 maps				
Agents	ODrM*	DCC	DHC	PRIMAL
4	50.00	48.575	52.33	79.08
8	52.17	59.60	63.90	76.53
16	59.78	71.34	79.63	107.14
32	67.39	93.54	100.10	155.21
64	82.60	135.55	147.26	170.48
Average steps in 80×80 maps				
Agents	ODrM*	DCC	DHC	PRIMAL
4	93.40	93.89	96.72	134.86
8	104.92	109.89	109.24	153.20
16	114.75	122.24	122.54	180.74
32	121.31	132.99	138.32	250.07
64	134.42	159.67	163.50	321.63
128	143.84	192.90	213.15	350.76

DHC are not included in this evaluation because they are all-to-all communication (more costly than request-reply).

A. Success Rate and Average Step

Success rate measures the ability to complete a MAPF task within given time steps. Average step measures the average time step consumed to finish a task, where smaller value indicates a better policy. We average both successful and unsuccessful cases to calculate the average steps. We set up two types of maps, 40×40 and 80×80 for testing. The obstacle density is set to be 0.3, the highest testing density in PRIMAL and DHC. We set up 200 test cases for each agent number in $\{4, 8, 16, 32, 64\}$. Additional agent number of 128 is tested in 80×80 maps due to a larger environment space. The maximum time step for 40×40 map is 256, and 386 for 80×80 map, the same as the PRIMAL's setting.

Fig. 3 shows the success rate of our method compared with PRIMAL and DHC in two types of environments. In general, DCC and DHC (RL plus communication) perform much better than PRIMAL (RL with expert guidance) in all cases. The performance of PRIMAL downgrades heavily in relatively larger environments (80×80) compared with smaller ones (40×40), which indicates the IL part of PRIMAL does not deliver good guidance to its RL component. Thus it suffers performance degradation on long-

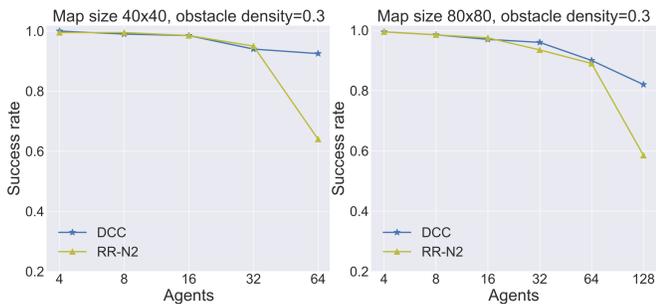


Fig. 4. Success rate of our method compared with RR-N2 in two different scenarios.

horizon task. In the environments where the agent density is low, such as 4, 8, 16 agents in both 40×40 and 80×80 maps, the difference between the success rate of DCC and DHC is small. In these environments, agents have a lower chance to meet and communicate with others, so the difference between broadcast communication (DHC) and selective communication (DCC) is not very significant. However, the success rate of DHC drops rapidly when agent density further grows. In agent dense environments, agents have a higher chance of encountering each other, but agents will obtain irrelevant and redundant information from others due to broadcast communication, which harms the cooperation among agents. By selecting and gathering relevant and influential information for communication, the cooperation is enhanced by DCC.

Table II verifies that DCC learns higher quality policies with respect to average steps. We include ODrM* [6], a centralized planner, as a reference. Especially compared with DHC, learning selective communication by DCC leads to much shorter paths in agent dense environments.

B. Selective Communication vs. Predefined Scope

Generally, selective communication by DCC does not need to set a communication scope. Any neighboring agent inside the FOV is possible to be requested by the central agent for communication. Previous broadcast communication method DHC defines the communication scope as the nearest two neighboring agents inside the FOV in order to reduce temporal redundant information. Although it is reasonable to assume the nearest agents are most relevant and influential, we show that selective communication outperforms the predefined scope (nearest two agents) for MAPF in both success rate and communication overhead.

To make a fair comparison, we develop a baseline model, named RR-N2 (request-reply, nearest two). The whole model architecture is similar to DCC. The only difference is that RR-N2 does not have the decision causal unit in Fig. 2, and in communication block, each agent only requests information from its two nearest agents inside the FOV. The training settings are the same as DCC as discussed in Section IV-C.

Fig. 4 shows the success rate of our method compared with this baseline. The difference in success rate is significant in agent-dense environments, i.e., 64 agents in 40×40 map and 128 agents in 80×80 map, which implies the nearest two agents are not always the most relevant ones. One may think

TABLE III
COMMUNICATION FREQUENCY IN TWO DIFFERENT ENVIRONMENTS
WITH OBSTACLE DENSITY = 0.3

Average Step Agents	Map size 40×40		Map size 80×80	
	DCC	RR-N2	DCC	RR-N2
4	2.42	36.88	1.06	18.36
8	11.56	209.79	5.75	105.86
16	60.47	959.38	24.98	469.58
32	294.69	4111.57	126.685	2125.94
64	1811.33	19490.09	562.11	8780.72
128	-	-	2915.84	36560.30

that DCC achieves a better performance due to unconstrained communication scope, where each agent may communicate with more than two agents inside its FOV and thus gathers more information than RR-N2. We show that this is not really the case and DCC actually learns succinct communication. In particular, we compute the average communication frequency in these test environments as shown in Table III. A pair of request and reply is counted as one time communication. The average communication overhead of RR-N2 is much more costly than DCC in all the testing cases, and is extremely expensive in agent-dense environments. By learning selective communication, DCC can greatly reduce communication by only focusing on relevant messages.

Combining these two results, we can conclude that although communication can help multi-agent cooperation for MAPF, agents should learn to actively select relevant agents for communication, because even the most nearest agents are not the relevant ones. Communicating with only two agents can already lead to huge communication overhead. Thus, researchers should always devise a succinct communication mechanism in order to reduce the communication overhead for easy deployment to real world problems.

Limitations: There is an additional computational cost associated with selecting neighboring agents. In extreme cases, the cost is of order of $O(n \times p)$. p is the FOV capacity.

VI. CONCLUSION

We propose DCC to learn succinct communication for MAPF with classical MAPF environment settings. Our aim is to enable agents to actively select relevant and influential neighbors for communication, instead of broadcasting. The selection is learned via the local policy adjustment effect between agents, which captures the necessity of communication. Empirical results show that selective communication with relevant agents improves the policy learning process. Moreover, DCC also serves as a component for communication reduction, greatly downscaling the communication overhead. Future work entails extensions of DCC to centralized training frameworks and cooperative Markov Game tasks. Another direction is to explore how to directly exclude sets of neighbors during neighboring selection to reduce complexity, and also the maximal number of agents that DCC can handle in different size of maps.

REFERENCES

- [1] J. Yu and S. M. LaValle, "Structure and intractability of optimal multi-robot path planning on graphs," in *Proc. AAAI Conference on Artificial Intelligence (AAAI'13)*, 2013.
- [2] J. Banfi, N. Basilico, and F. Amigoni, "Intractability of time-optimal multirobot path planning on 2d grid graphs with holes," *IEEE Robotics and Automation Letters*, vol. 2, no. 4, pp. 1941–1947, 2017.
- [3] J. Yu and S. M. LaValle, "Planning optimal paths for multiple robots on graphs," in *Proc. International Conference on Robotics and Automation (ICRA'13)*. IEEE, 2013, pp. 3612–3617.
- [4] P. Surynek, A. Felner, R. Stern, and E. Boyarski, "Efficient sat approach to multi-agent path finding under the sum of costs objective," in *Proc. European Conference on Artificial Intelligence (ECAI'16)*, 2016, pp. 810–818.
- [5] M. Goldenberg, A. Felner, R. Stern, G. Sharon, N. Sturtevant, R. C. Holte, and J. Schaeffer, "Enhanced partial expansion a*," *Journal of Artificial Intelligence Research*, vol. 50, pp. 141–187, 2014.
- [6] G. Wagner and H. Choset, "Subdimensional expansion for multirobot path planning," *Artificial Intelligence*, vol. 219, pp. 1–24, 2015.
- [7] G. Sharon, R. Stern, A. Felner, and N. R. Sturtevant, "Conflict-based search for optimal multi-agent pathfinding," *Artificial Intelligence*, vol. 219, pp. 40–66, 2015.
- [8] J. Li, A. Felner, E. Boyarski, H. Ma, and S. Koenig, "Improved heuristics for multi-agent path finding with conflict-based search," in *Proc. International Joint Conference on Artificial Intelligence (IJCAI'19)*, vol. 2019, 2019, pp. 442–449.
- [9] G. Sartoretto, J. Kerr, Y. Shi, G. Wagner, T. S. Kumar, S. Koenig, and H. Choset, "Primal: Pathfinding via reinforcement and imitation multi-agent learning," *IEEE Robotics and Automation Letters (RA-L)*, vol. 4, no. 3, pp. 2378–2385, 2019.
- [10] B. Wang, Z. Liu, Q. Li, and A. Prorok, "Mobile robot path planning in dynamic environments through globally guided reinforcement learning," *IEEE Robotics and Automation Letters (RA-L)*, vol. 5, no. 4, pp. 6932–6939, 2020.
- [11] B. Riviere, W. Hönig, Y. Yue, and S.-J. Chung, "Glas: Global-to-local safe autonomy synthesis for multi-robot motion planning with end-to-end learning," *IEEE Robotics and Automation Letters (RA-L)*, vol. 5, no. 3, pp. 4249–4256, 2020.
- [12] Z. Liu, B. Chen, H. Zhou, G. Koushik, M. Hebert, and D. Zhao, "Mapper: Multi-agent path planning with evolutionary reinforcement learning in mixed dynamic environments," in *Proc. IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'20)*. IEEE, 2020, pp. 11 748–11 754.
- [13] Q. Li, F. Gama, A. Ribeiro, and A. Prorok, "Graph neural networks for decentralized multi-robot path planning," in *Proc. IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'20)*. IEEE, 2020, pp. 11 785–11 792.
- [14] Q. Li, W. Lin, Z. Liu, and A. Prorok, "Message-aware graph attention networks for large-scale multi-robot path planning," *IEEE Robotics and Automation Letters (RA-L)*, vol. 6, no. 3, pp. 5533–5540, 2021.
- [15] Z. Ma, Y. Luo, and H. Ma, "Distributed heuristic multi-agent path finding with communication," in *Proc. IEEE International Conference on Robotics and Automation (ICRA'21)*, 2021.
- [16] M. L. Littman, "Markov games as a framework for multi-agent reinforcement learning," in *Proc. International Conference on Machine Learning (ICML'94)*, 1994, pp. 157–163.
- [17] A. Singh, T. Jain, and S. Sukhbaatar, "Learning when to communicate at scale in multiagent cooperative and competitive tasks," in *Proc. International Conference on Learning Representations (ICLR'19)*, 2019.
- [18] D. Kim, S. Moon, D. Hostallero, W. J. Kang, T. Lee, K. Son, and Y. Yi, "Learning to schedule communication in multi-agent reinforcement learning," in *Proc. International Conference on Learning Representations (ICLR'19)*, 2019.
- [19] S. Q. Zhang, Q. Zhang, and J. Lin, "Efficient communication in multi-agent reinforcement learning via variance based control," in *Proc. Advances in Neural Information Processing Systems (NeurIPS'19)*, 2019, pp. 3230–3239.
- [20] S. Q. Zhang, Q. Zhang, and J. Lin, "Succinct and robust multi-agent communication with temporal message control," in *Proc. Advances in Neural Information Processing Systems (NeurIPS'20)*, 2020.
- [21] T. Wang, J. Wang, C. Zheng, and C. Zhang, "Learning nearly decomposable value functions via communication minimization," in *Proc. International Conference on Learning Representations (ICLR'20)*, 2020.
- [22] Z. Ding, T. Huang, and Z. Lu, "Learning individually inferred communication for multi-agent cooperation," in *Proc. Advances in Neural Information Processing Systems (NeurIPS'20)*, 2020.
- [23] M. Everett, Y. F. Chen, and J. P. How, "Motion planning among dynamic, decision-making agents with deep reinforcement learning," in *Proc. IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'18)*. IEEE, 2018, pp. 3052–3059.
- [24] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *Proc. International Conference on Machine Learning (ICML'16)*. PMLR, 2016, pp. 1928–1937.
- [25] A. Das, T. Gervet, J. Romoff, D. Batra, D. Parikh, M. Rabbat, and J. Pineau, "Tarmac: Targeted multi-agent communication," in *Proc. International Conference on Machine Learning (ICML'19)*. PMLR, 2019, pp. 1538–1546.
- [26] J. Jiang, C. Dun, T. Huang, and Z. Lu, "Graph convolutional reinforcement learning," in *Proc. International Conference on Learning Representations (ICLR'20)*, 2020.
- [27] Y. Niu, R. Paleja, and M. Gombolay, "Multi-agent graph-attention communication and teaming," in *Proc. International Conference on Autonomous Agents and MultiAgent Systems (AAMAS'21)*, 2021.
- [28] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. Advances in neural information processing systems (NeurIPS'17)*, 2017, pp. 5998–6008.
- [29] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE transactions on neural networks*, vol. 20, no. 1, pp. 61–80, 2008.
- [30] M. Barer, G. Sharon, R. Stern, and A. Felner, "Suboptimal variants of the conflict-based search algorithm for the multi-agent pathfinding problem," in *Annual Symposium on Combinatorial Search*, 2014.
- [31] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson, "Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning," in *Proc. International Conference on Machine Learning (ICML'18)*. PMLR, 2018, pp. 4295–4304.
- [32] R. Stern, N. R. Sturtevant, A. Felner, S. Koenig, H. Ma, T. T. Walker, J. Li, D. Atzmon, L. Cohen, T. S. Kumar, et al., "Multi-agent pathfinding: Definitions, variants, and benchmarks," in *Twelfth Annual Symposium on Combinatorial Search*, 2019.
- [33] M. Damani, Z. Luo, E. Wenzel, and G. Sartoretto, "Primal₂: Pathfinding via reinforcement and imitation multi-agent learning-lifelong," *IEEE Robotics and Automation Letters (RA-L)*, vol. 6, no. 2, 2021.
- [34] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *Proc. International Conference on Machine Learning (ICML'93)*, 1993, pp. 330–337.
- [35] P. Sunehag, G. Lever, A. Grusl, W. M. Czarnecki, V. F. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, and T. Graepel, "Value-decomposition networks for cooperative multi-agent learning based on team reward," in *Proc. International Conference on Autonomous Agents and MultiAgent Systems (AAMAS'18)*, 2018, pp. 2085–2087.
- [36] S. Kullback, *Information theory and statistics*. Courier Corporation, 1997.
- [37] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [38] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Dueling network architectures for deep reinforcement learning," in *Proc. International Conference on Machine Learning (ICML'16)*. PMLR, 2016, pp. 1995–2003.
- [39] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proc. International Conference on Machine Learning (ICML'09)*, 2009, pp. 41–48.
- [40] S. Kapturovski, G. Ostrovski, J. Quan, R. Munos, and W. Dabney, "Recurrent experience replay in distributed reinforcement learning," in *Proc. International Conference on Learning Representations (ICLR'18)*, 2018.
- [41] D. Horgan, J. Quan, D. Budden, G. Barth-Maron, M. Hessel, H. van Hasselt, and D. Silver, "Distributed prioritized experience replay," in *Proc. International Conference on Learning Representations (ICLR'18)*, 2018.
- [42] P. Moritz, R. Nishihara, S. Wang, A. Tumanov, R. Liaw, E. Liang, M. Elibol, Z. Yang, W. Paul, M. I. Jordan, et al., "Ray: A distributed framework for emerging AI applications," in *USENIX Symposium on Operating Systems Design and Implementation*, 2018, pp. 561–577.