

---

# MAIN: A Real-world Multi-agent Indoor Navigation Benchmark for Cooperative Learning

---

Fengda Zhu<sup>1</sup> Siyi Hu<sup>1</sup> Yi Zhang<sup>1</sup> Haodong Hong<sup>1</sup>  
Yi Zhu<sup>3</sup> Xiaojun Chang<sup>1</sup> Xiaodan Liang<sup>2,4</sup>  
<sup>1</sup>Monash University <sup>2</sup> Sun Yat-sen University  
<sup>3</sup> Huawei Noah’s Ark Lab <sup>4</sup> Dark Matter AI Inc.

## Abstract

1 The ability to cooperate and work as a team is one of the ‘holy grail’ goals of  
2 intelligent robots. Previous works have proposed many multi-agent reinforcement  
3 learning methods to study this problem in diverse multi-agent environments. How-  
4 ever, these environments have two limitations, which make them unsuitable for  
5 real-world applications: 1) the agent observes clean and formatted data from the  
6 environment instead of perceiving the noisy observation by themselves from the  
7 first-person perspective; 2) large domain gap between the environment and the  
8 real world scenarios. In this paper, we propose a Multi-Agent Indoor Navigation  
9 (MAIN) benchmark<sup>1</sup>, where agents navigate to reach goals in a 3D indoor room  
10 with realistic observation inputs. In the MAIN environment, each agent observes  
11 only a small part of a room via an embodied view. Less information is shared  
12 between their observations and the observations have large variance. Therefore,  
13 the agents must learn to cooperate with each other in exploration and communi-  
14 cation to achieve accurate and efficient navigation. We collect a large-scale and  
15 challenging dataset to research on the MAIN benchmark. We examine various  
16 multi-agent methods based on current research works on our dataset. However,  
17 we find that the performances of current MARL methods does not improve by the  
18 increase of the agent amount. We find that communication is the key to addressing  
19 this complex real-world cooperative task. By Experimenting on four variants of  
20 communication models, we show that the model with recurrent communication  
21 mechanism achieves the best performance in solving MAIN.

## 22 1 Introduction

23 Cooperative multi-agent problems are ubiquitous in real-world applications, for example, multiplayer  
24 games [40, 38, 18], multi-robot control [29], language communication [48, 15, 33], and social  
25 dilemmas [23]. These applications focus on solving the sequential decision-making problem of  
26 multiple autonomous agents within a common environment, which could be systematically modeled  
27 as the multi-agent reinforcement learning (MARL) paradigm [34, 52, 65]. Compared to traditional  
28 reinforcement learning, MARL has two major challenges. The first is the partial observability. Each  
29 agent observes only part of the global state. The second is the instability of learning decentralised  
30 policies. Recent works have proposed diverse environments to validate the effectiveness of the  
31 MARL algorithms, such as grounded communication environment [33], StarCraft II [42], DOTA2 [7],  
32 multi-agent emergence environments [4], soccer shooting [26], etc.

33 Most of these game-based MARL environments are quite different from the real-world situation such  
34 as robotics and auto-driving. For example, the agent observes clean and formatted data from the

---

<sup>1</sup><http://main-dataset.github.io/>

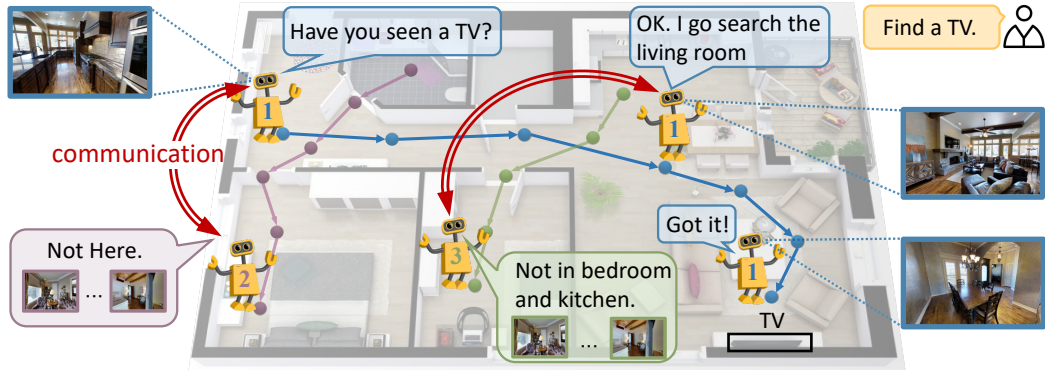


Figure 1: A demonstration of our Multi-agent Indoor Navigation (MAIN) benchmark. The agent 1 moves in the room to find the target by cooperating and communicating with agent 2 and agent 3.

35 game environment instead of perceiving realistic observations by themselves from the first-person  
 36 perspective. The transition function in the game environments is based on simple rules without  
 37 simulating the physical rules and considering the interference in the real world. Therefore, there is  
 38 large domain gap between current MARL environments and real world, which limits current MARL  
 39 models [27, 40, 62] to be applied to real-world scenarios.

40 *There are no good models without good data* [41]. To overcome these limitations, we propose a novel  
 41 benchmark, Multi-Agent Indoor Navigation (MAIN), where multiple agents are required to navigate  
 42 to reach goals in a 3D indoor room. To obtain realistic observations from the first-person view of  
 43 agents, we adopt Habitat [44] simulator to render realistic egocentric RGB-D image observations  
 44 for agents. At each navigation step, each agent observes an RGB-D image from its own first-person  
 45 perspective and makes action decision including, ‘turn left’, ‘turn right’ and ‘step forward’. The  
 46 setting of realistic egocentric observation is closer to real-world situations, which makes the learned  
 47 agents easier to be transferred to real-world applications like robotics.

48 In other MARL environments such as StarCraft II [42] and multi-agent emergence environments [4],  
 49 the observations of different agents have a large proportion of overlapping. For example, the status  
 50 (position or health) of an agent or an object is fully observed if it is located within the vision  
 51 range of another agent. It is unrealistic compared to the real-world where the status of agents  
 52 and objects is not fully observable. In our proposed MAIN environment, the appearance, shape,  
 53 and size of an object will be very different, when observed by the agents from different angles,  
 54 especially in the first-person views where the angles are dynamically changing all the time. In  
 55 addition to the realistic observation, MAIN adopts the Bullet physics engine [12] to provide a more  
 56 realistic transition function. Unlike other environments [42, 4, 26], where the agent receives the  
 57 high-precision localization information from the environment, MAIN does not provide a compass  
 58 sensor and requires the agent to navigate solely using an egocentric RGB-D camera. Compared with  
 59 previous single-agent navigation environments such as MINOS [43] and Habitat [44], we implement  
 60 a asynchronous-synchronous pipeline for efficient multi-agent data sampling.

61 To the best of our knowledge, our MAIN is the first multi-agent real-world navigation environment.  
 62 The environments of MAIN benchmark bring new challenges, such as learning from realistic obser-  
 63 vations and less observation overlapping between agents. These new challenges raise additional  
 64 requirement for better utilizing the information by sharing the individual observation with other  
 65 collaborators and making decisions based on both self-observation and collaborators’ observation.  
 66 However, many MARL methods [40, 1, 18, 55] adopt the Centralized Training Decentralized Ex-  
 67 ecution (CTDE) [36] framework, which forbids the real-time information sharing among agents.  
 68 Therefore, these methods are not suitable for a real-world simulated environment like MAIN. To  
 69 address this, we propose a new cooperative multi-agent communication mechanism to enable the  
 70 agents to exchange the information in a real-time manner. This communication mechanism may not  
 71 be essential for game-based or highly-simplified tasks like Hanabi [5], SMAC [53] and Hide-and-  
 72 seek [4] but is crucial in real-world tasks like MAIN due to the highly unlinear and little overlapped  
 73 observations. In short, being aware of the status of the collaborators is beneficial to making a good  
 74 action decision, and is critical for accomplishing a real-world cooperative task. An overview of  
 75 MAIN task with the communication mechanism is shown in Fig. 1. Three agents start from different

Input	Environment	Multi-agent	Embodied	Observation Overlap	Physics Engine
Array	Hanabi [5]	✓	-	$\frac{N-1}{N}$	-
	Diplomacy [37]	✓	-	$\frac{N-1}{N}$	-
	EGCL [33]	✓	✗	100%	-
	DOTA2 [7]	✓	✗	100%	Rubikon
	SMAC [42]	✓	✗	81.0%-89.8%	Havok
	Hide-and-seek [4]	✓	✗	$\frac{1}{N}$ -100%	MuJoCo
	Soccer [26]	✓	✗	100%	MuJoCo
Syn image	MARLÖ [39]	✓	✓	$\frac{1}{N}$	hybrid voxel
	AI2-THOR [21]	✗	✓	$\frac{1}{N}$	Unity
Real image	Habitat [44]	✗	✓	$\frac{1}{N}$	Bullet
	MAIN (Ours)	✓	✓	$\frac{1}{N}$	Bullet

Table 1: Compared with existing MARL environments (N is the number of agents). The egocentric view means if the agent has a local observation of environment from its perspective rather than receiving global information. The embodied view represents whether an agent observe the 3D environment from the first-person perspective (we only compare the navigable environments).

76 positions and are asked to find a TV in this room. Each agent receives a first-person photo-realistic  
77 observation from their perspective respectively. They explore the room and communicate with each  
78 other to exchange their discoveries. By active exploration and cooperative communication, the No. 1  
79 agent finally finds the TV.

80 Considering that agents can have different targets at the same time, we propose two sub-tasks for  
81 our MAIN task: 1) **Shared-target navigation** where all agents are asked to find a shared target  
82  $g$ ; 2) **Individual-target navigation** where each agent  $i$  has its own target  $g_i$ . In the shared-target  
83 navigation sub-task, we mainly evaluate the agents cooperation ability of searching separately for  
84 a target. In the individual-target navigation, we focus on evaluating the ability of cooperative  
85 information exchanging. To fully investigate the MAIN benchmark, we construct a large-scale  
86 dataset consists of 24M episodes within 90 houses for training, validation, and testing splits, which is  
87 10 times larger than the dataset in [44]. The data is automatically labeled within the environment.  
88 Compared with other datasets [57, 44], our dataset is challenging since it provides more long-term  
89 hard samples.

90 Along side with the environment, we provide multiple baselines for MARL research community for  
91 fast evaluating their effectiveness for real-world deployment. We build our baseline models based on  
92 previous MARL works [27, 1, 62] to validate the effectiveness of our benchmark and dataset. We  
93 find that the number of agents improves the navigation performance in simple baselines. However,  
94 without communication, the navigation performance will not increase by increasing the number of  
95 agents. And we find it is essential for agents to communicate with each other in addressing complex  
96 real-world cooperation tasks. We experiment on four kinds of communication variants and find that  
97 the model with a recurrent actor-critic mechanism to encode historical communication messages  
98 significantly outperforms other models. In summary, we make the following contributions: 1) we  
99 propose the MAIN benchmark to research on multi-agent problem in a realistic environment; 2) we  
100 collect a large-scale and challenging dataset and benchmark several MARL baseline models; 3) we  
101 propose a communication module that benefits for the real-world multi-agent system.

## 102 2 Related Work

103 **Multi-agent Environments** have been proposed to research multi-agent problems. However, previ-  
104 ous works ignore the importance of implementing a realistic environment, which limits the learned  
105 model to be applied on real-world applications such as robotics. In Tab. 1, we compare the differences  
106 between our MAIN environment with previous multi-agent environments. To the best of our knowl-  
107 edge, we claim that our MAIN environment is the first multi-agent environment that offers realistic  
108 image input. Previous works [7, 42, 4] get clean and formatted array data via programming interfaces.  
109 Therefore, it would be hard for the learned model to overcome the challenge of the interference  
110 from noisy data. MARLÖ [39] provide synthetic image whose domain is largely deviated from the  
111 scenarios of the real-world.

112 Our model provides an egocentric and embodied view, which is quite applicable for robots in the  
113 real world. A particular challenge in the MARL problem is partial observability. Each agent could  
114 only observe a part of the global state and solve the problem by communication. However, our  
115 investigation reveals that some of the previous MARL environments [42, 4], even though claimed  
116 to be partially observable, have large observation overlap. Little observations overlap makes our  
117 benchmark more challenging than the previous environments, because little information sharing  
118 makes cooperation difficult. A realistic physical engine helps simulate a real-world transition function.  
119 Our environment adopt Bullet engine to simulate physics activities such as acceleration and collision.

120 **Multi-agent Reinforcement Learning** extends the problem of reinforcement learning [32, 20, 50]  
121 into multi-agent scenario and brings new challenges. Most MARL problems [28, 47, 65] falls into the  
122 centralized training with decentralized execution (CTDE) architecture [36, 22, 61]. Some works are  
123 dedicated to improving the mixing network which mixes the agent network outputs to learn a joint  
124 action-value function [40, 49, 55]. Other works are aiming at improving the network structure and  
125 developing a individual function with better representation and transfer capability [16, 18]. Besides,  
126 the utilization of state information varies. IPPO [13] incorporates the global state information barely  
127 by sharing network parameters among critics of individual agents. While MAPPO [63] constructs  
128 a centralised value function upon agents which takes the aggregated global state information as  
129 inputs. Researchers find that communication is critical cooperative multi-agent problems. Lowe  
130 *et al.* [27] propose MADDPG, a framework based on DDPG [24] with cooperative value function.  
131 Later, R-MADDPG [54] equips with recurrent actor critic models, simultaneously learning policies for  
132 navigation and communication towards better information utilization and resource distribution. We  
133 implement diverse methods to support extensive research and illustrate the novelties and challenges  
134 of MAIN.

135 **Embodied Navigation Environments.** Simulations such as Matterport3D simulator [3], Gibson  
136 simulator [60] and Habitat [44] propose high-resolution photo-realistic panoramic view to simulate  
137 more realistic environment. Rendering frame rate is also important to embodied simulators since it is  
138 critical to training efficiency. MINOS [43] runs more than 100 frame per second (FPS), which is 10  
139 times faster than its previous works. Habitat [44] runs more than 1000 FPS on  $512 \times 512$  RGB+depth  
140 image, making it become the fastest simulator among existing simulators. Some complex tasks may  
141 require a robot to interact with objects, such as picking up a cup, moving a chair or opening a door.  
142 AI2-THOR [21], iGibson [59] and RoboTHOR [14] provide interactive environments to train such a  
143 skill. Multi-agent reinforcement learning [25, 51] is a rising problem of cooperation and competition  
144 among agents. Based on the Habitat simulator, we construct a multi-agent environment to research  
145 on realistic MARL problem.

146 **Embodied Navigation Learning** is attracting rising attention in the community and lots of methods  
147 have been proposed to address this problem. Based on conventional reinforcement learning meth-  
148 ods [31], Wu *et al.* [58] introduce an LSTM layer to encode the historical information. Wang *et al.*  
149 [56] propose to jointly learn a navigation model with imitation learning and supervised learning.  
150 Some works [19, 30, 66] propose auxiliary tasks to exploit extra training signals for learning naviga-  
151 tion. SLAM-based methods [64, 9, 10] are widely adopted in navigation due to its capability of  
152 modeling the room structure. Nonetheless, those tasks do not conducted in multi-agent setting which  
153 requires cooperation and communication, and consequently, being more flexible and practical.

## 154 3 Multi-agent Indoor Navigation Benchmark

### 155 3.1 Task Definition

156 Here we define our proposed Multi-agent Indoor Navigation (MAIN) Benchmark in detail. MAIN  
157 requires multiple agents  $E = \{e_1, \dots, e_n\}$  to navigate to reach a set of targets accurately and efficiently  
158 in an indoor environment. At the beginning of an episode, each agent is told to reach a target  $g_i$ . For  
159 each step, the agent observes an observation and make an action decision. The observation contains  
160 an RGB-D image, localization information from a GPS compass and contact information from a  
161 physics sensor. An action could be ‘turn left’, ‘turn right’, ‘step forward’ and ‘found’. The agent uses  
162 the first three actions to navigate in the environment and uses the last action to declare it has found the  
163 target object. The episode of this agent is considered succeed if the ‘found’ action is selected while  
164 the agent is located with the threshold toward the target object. Otherwise, the episode is consider a

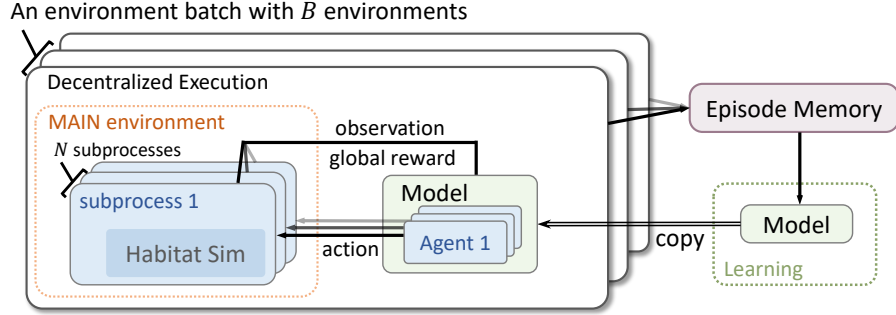


Figure 2: Training overview of our MAIN environment.

165 failure if the ‘found’ action is slected while the agent is out of range or it has navigated for maximum  
 166 steps without finding the target.

167 Based on the above rules, we propose two sub-tasks, shared-target navigation and individual-target  
 168 navigation. **Shared-target navigation** is a task where all agents are asked to find a shared target  $g$ .  
 169 The task is considered succeed if any agent reaches  $g$  and considered failure if any agent fails within  
 170 its episode. **Individual-target navigation** is a task where each agent  $i$  has its own target  $g_i$ , and the  
 171 task is considered succeed only when all agents successfully find its own target. We do experiments  
 172 under the setting of different agents and different tasks to demonstrate the challenge and novelty of  
 173 this benchmark. The reward function is defined by shortened distance similar to [57].

### 174 3.2 Multi-agent Indoor Navigation Environment

175 This environment is built based on Habitat [44] simulator. Habitat simulator renders the 3D assets  
 176 of an house and provide a photo-realistic embodied environment for agents. The Habitat simulator  
 177 provide multiple sensors including RGB-D image, GPS compass and contact. The Habitat is built  
 178 upon the Bullet physics engine that enables realistic graphics rendering, velocity and acceleration  
 179 simulation, and contact simulation. However, the rendering process is computation consuming and  
 180 time costly. Therefore, we design a asynchronous-synchronous pipeline for data efficiency.

181 Our pipeline is shown in Fig. 2. The MAIN environment creates  $B$  sub-environments for decentralized  
 182 execution to sample data for training, where  $B$  is the size of the minibatch. Each sub-environment  
 183 creates  $N$  processes, where the  $N$  is the number of agents. Each process has a copy of a Habitat  
 184 simulator, and each simulator individually simulates the state of an agent and renders the RGB-D  
 185 image observation for an agent. The MAIN sub-environment synchronizes the processes and interacts  
 186 with a copy of a multi-agent navigation model. In the decentralized execution, the parameters of the  
 187 the multi-agent navigation model are shared across all sub-environments. The multi-agent navigation  
 188 model predicts actions for each agent for each step. The predicted actions are sent to the MAIN  
 189 environment and then distributed to each process to execute. The Habitat simulator execute the action  
 190 and return the updated state and the current partial observation to the MAIN sub-environment. The  
 191 MAIN sub-environment calculate the global reward based on the global state and send the global  
 192 reward and observations to the model. For each step, the global reward, observations for all agents  
 193 and actions that agents predict are stored in the episode memory. We sample the episodes from the  
 194 episode memory to optimize a centralized model by stochastic gradient descent (SGD). The model  
 195 after a step of SGD optimization is copied to each MAIN sub-environment to update the execution  
 196 model.

### 197 3.3 Data Collection

198 We use the room textures and other 3D assets provided by Matterport3D [8] to build the MAIN  
 199 environment. Matterport3D consists of 10,800 panoramic views constructed from 194,400 RGB-D  
 200 images of 90 building-scale scenes, where 61 scenes for training, 11 for validation, and 18 for testing  
 201 following the standard split [8]. We provide episode data for learning and testing. An episode is  
 202 defined by a starting position where the agent starts and the target position where the agent is required  
 203 to reach. Both the starting position and the target positions are randomly sampled from navigable  
 204 points within an environment. We ensures there is at least one navigable path from the starting

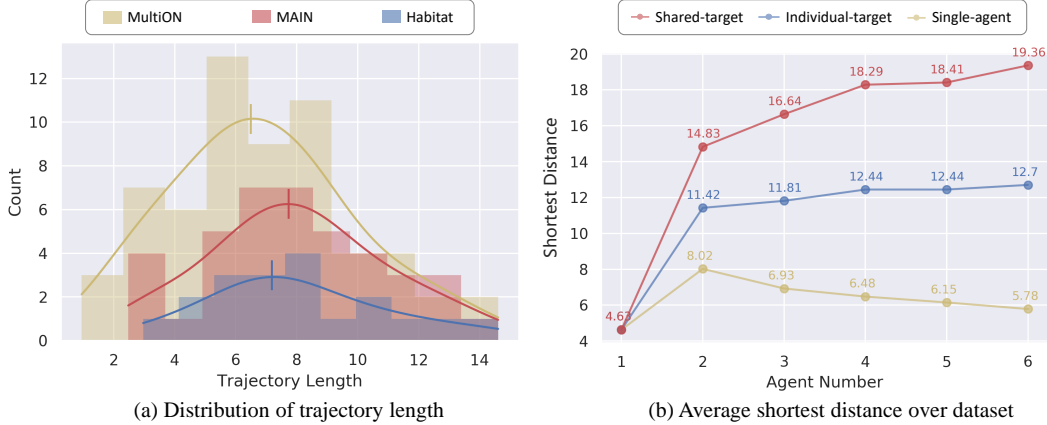


Figure 3: An analysis of our MAIN dataset. The curve in the left figure stands for the Gaussian smoothing. The lines are the mean values of the smoothed Gaussian distributions.

205 position to the target. And we constrain the length of an episode to be between 2m and 20m, which  
 206 ensures that each episode is neither too trivial nor too hard.

207 We compare the distribution of the average trajectory length within a room with MultiION [57] and  
 208 the object navigation data in the Habitat Challenge [6] in Fig. 3(a). Due to the different structures  
 209 of the house, the episode data from each house have different average length. We find that with the  
 210 same room setting, our average trajectory length is longer than both MultiION and Habitat, proving  
 211 that our data is more challenging. Our dataset provide 24M episodes, 10 times more than the data  
 212 scale of the Habitat dataset.

213 The Fig. 3(b) shows the average distance that the agents need to navigate to successfully accomplish  
 214 task. It reveals the gap of the difficulty among the two sub-tasks and the single-agent navigation task  
 215 accompany with the agent amount using our dataset. With the increase of the agent amounts, the  
 216 difficulty of individual-target task is significantly increasing while the shared-target task is reducing.

### 217 3.4 Metrics

218 The MAIN task is evaluated from two aspects: navigation accuracy and efficiency. We use the  
 219 following metrics to quantitatively measure the effectiveness of models:

220 **Success Rate** is used to measure if the agent successfully finds the target when it yields ‘found’. The  
 221 agent is regarded ‘success’ only if it is located within a threshold distance towards the target.

222 **Distance** indicates the average distance forward the target when the agent stops. This metric is useful  
 223 when the success rate is low.

224 **SPL**, short for Success weighted by Path Length [2], evaluates the accuracy and efficiency simultane-  
 225 ously. The SPL is calculated by  $\frac{1}{N} \sum_{i=1}^N S_i \frac{p_i}{l_i}$ , where the  $N$  is all testing samples,  $S_i$  is the success  
 226 indicator,  $p_i$  is the shortest path length, and the  $l_i$  is the actual path length in testing. We adopt the  
 227 SPL as our main metric.

## 228 4 Multi-agent Models

### 229 4.1 Preliminaries

230 We systematically model our MAIN problem as a multi-agent reinforcement learning paradigm  
 231 which is described as a partially-observed Markov decision process (POMDP) [35].  $P(s'|s, a)$  is the  
 232 transition probability that transforms the current state space  $S$  to the next state space  $S'$  conditions on  
 233 the a global action  $a \in A$ . We follow the centralized-training decentralized-execution framework  
 234 that parameterize the shared policy of each agent as  $\pi_\theta$ . For each step  $t$ , the agent  $i$  receive its partial  
 235 observation  $o_{t,i}$  and choose its action by  $a_{t,i} = \pi_{\theta_i}(o_{t,i})$ . The global action  $a_t = a_{1,t}, \dots, a_{n,t}$ . All  
 236 agents share the same global reward function  $r(s, a) : S \times A \rightarrow \mathbb{R}$ . And the  $\gamma \in [0, 1)$  is a discount

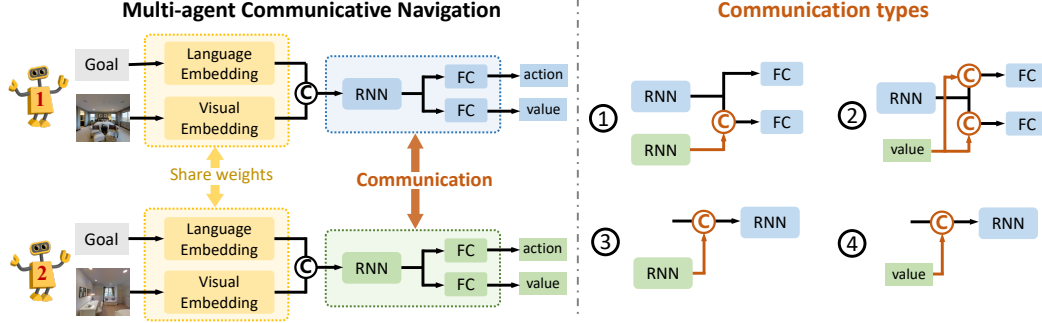


Figure 4: A demonstration of our multi-agent communicative navigation framework. We take a two-agent model for instance. A dot bounding box denote an agent. The dual-directed orange arrow stands for weight sharing between the orange boxes. The orange arrow stands for communication. On the right side, we show four kinds of communication variants.

237 factor that defines the length of the horizon. We optimization the parameter  $\theta$  by minimizing the  
 238 optimization objective  $J(\theta) = \mathbb{E} [\sum_t \gamma^t r(s_t, a_t)]$  by PPO algorithm [46].

## 239 4.2 Baseline Multi-agent Models

240 We implement several multi-agent models to investigate the performance of multi-agent models on  
 241 the MAIN task.

242 **Random navigator with oracle founder.** We implement a random baseline which randomly sample  
 243 the action of ‘turn left’, ‘turn right’ and ‘go forward’. And the baseline model has the oracle ‘found’  
 244 module that yields ‘found’ as long as the agent reaches within the success range of navigation. This  
 245 baseline model is used to validate if our dataset is too easy or have severe bias.

246 **Multi-single agent.** This model is implement in the PPO [46] that is trained in a single-agent  
 247 paradigm but tested in a multi-agent paradigm. We research on this model to see if the number of  
 248 agents help the performance of navigaion in multi-agent paradigm.

249 **IPPO.** The IPPO model learns the global reward and share network parameters each agent. The  
 250 difference between PPO and IPPO [1] is that the PPO model receives a single-agent reward while the  
 251 IPPO model receives a global reward that influenced by other agents. The actions of other agents  
 252 cause the instability of the global reward, which increases the difficulty of training.

253 **MAPPO.** Based on IPPO, MAPPO [63] introduces a centralised value function upon agents with  
 254 global state inputs. However, the original MAPPO does not consider the importance of encoding the  
 255 historical communicative information, which limits its application in complex environments where  
 256 the observations of the agents have little in common and the historical information is important in  
 257 action decision. In our implementation, the CNN and RNN are shared among agents while the each  
 258 agent has its own actor and critic functions.

## 259 4.3 Multi-agent Cooperative Communication Navigation

260 In this section, we are going to introduce our cooperative communicative navigation model, as shown  
 261 in Fig. 4. We take a two-agent situation for demonstration. The framework firstly embeds the target as  
 262 an embedding feature, and extracts visual feature using an Convolutional Neural Network (CNN) [17]  
 263 module. The parameters of the embedding layer and the CNN layer are shared between agents to  
 264 ensure the generazability. Then the target feature and visual feature are concatenated to feed the  
 265 Recurrent Neural Network (RNN) [11] module. The RNN module is adopted to encode historical  
 266 information. Since the agents receive partial observation, it is important to memorize the previous  
 267 observation to help the agent build a more comprehensive understanding of the environment. The  
 268 historical feature from the RNN is send to two fully connected layers. One outputs an probability  
 269 that represents the preference of making action decision and the other predicts a value to estimate  
 270 the effective of the current situation. The model is optimized by PPO algorithm. To be specific, the  
 271 action prediction is supervised by policy gradient loss and the value prediction is supervised by the  
 272 bellman equation.

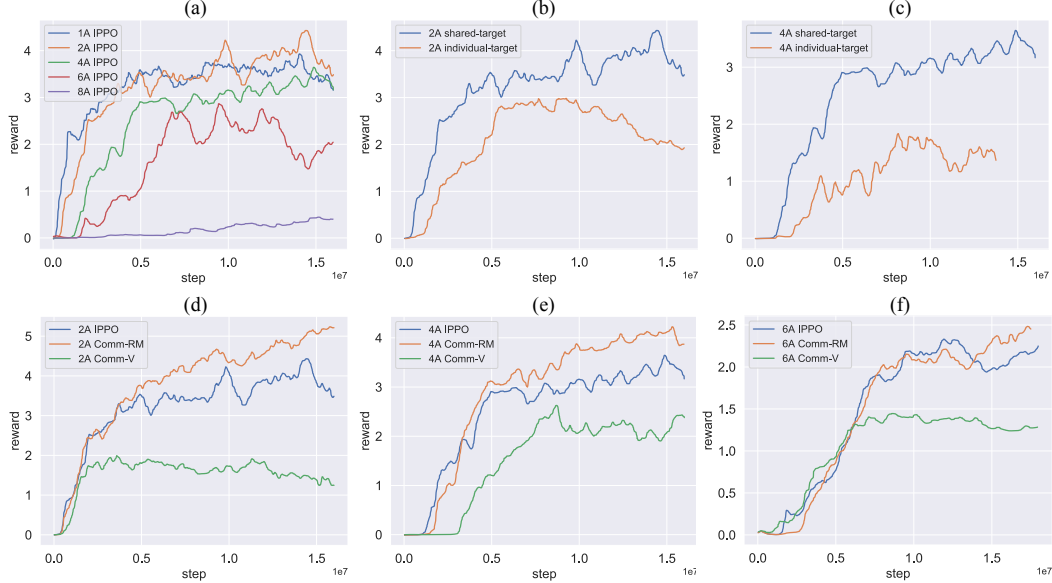


Figure 5: The result curves of our experiments.

Models	2 Agents				3 Agents			
	Length	Distance	Success rate	SPL	Length	Distance	Success rate	SPL
Random	3.39	12.75	0.00	0.00	3.30	16.67	0.00	0.00
Multi-PPO [46]	232.89	10.66	0.21	0.17	47.11	16.10	0.01	0.00
IPPPO [1]	256.67	15.55	0.08	0.06	137.24	15.94	0.02	0.01
Comm-S	351.03	<b>10.02</b>	0.12	0.06	75.92	16.16	0.05	0.05
Comm-V	68.14	12.02	0.03	0.03	80.23	<b>10.36</b>	0.05	0.04
Comm-RM	309.7	12.78	<b>0.3</b>	<b>0.23</b>	312.86	14.59	<b>0.13</b>	<b>0.06</b>
Comm-RV	298.1	11.56	0.24	0.17	301.2	12.32	0.08	0.05

Table 2: The testing results of different models. Multi-PPO: single-agent PPO model tested in multi-agent environment. The four variants of our communicative models is denoted as Sequential Communication model (Comm-S), Value Communication model (Comm-V), recurrent message communication model (Comm-RM), and recurrent value communication model (Comm-RV).

273 The agents exchange information between the the blue block and the green block to obtain more  
 274 knowledge and build a more comprehensive understanding of the environment. The feature vector that  
 275 an agent send is named as ‘message’. The agent that receives the message is named the ‘receiver’ and  
 276 the agent that sends the message is named the ‘sender’. The gradient is not back-propagated from the  
 277 ‘receiver’ to the ‘sender’ since it causes severe instability in training, which makes the performance of  
 278 the learned navigation model to be almost zero. On the left we show four communicative variants. We  
 279 name the them as sequential communication model (Comm-S) , value communication model (Comm-  
 280 V), recurrent message communication model (Comm-RM), and recurrent value communication  
 281 model (Comm-RV).

## 282 5 Experiment

283 **Implementation Details** Our communicative model is built based on our implementation of [1].  
 284 We train all of our models for 15M iterations. We adopt Adam optimizer whose learning rate is  
 285  $2.5 \times 10^{-4}$ . The discount factor  $\gamma = 0.99$  and the TD( $\lambda$ ) factor in GAE [45] is 0.95. Our model is  
 286 trained on by 8 GPUs (7 GPUs for rendering image inputs and 1 GPU for optimization) for 36 hours.

287 **Ablation for Agent Amount** The Fig. 5(a) ablates the amount of agent in MARL learning. We find  
 288 that with the amount increasing, the navigation performance is declining. More agent narrows the  
 289 searching area for find a target. However, the global reward is easily effected by the actions of other  
 290 agents, and therefore, hard to give an agent a clear guidance.



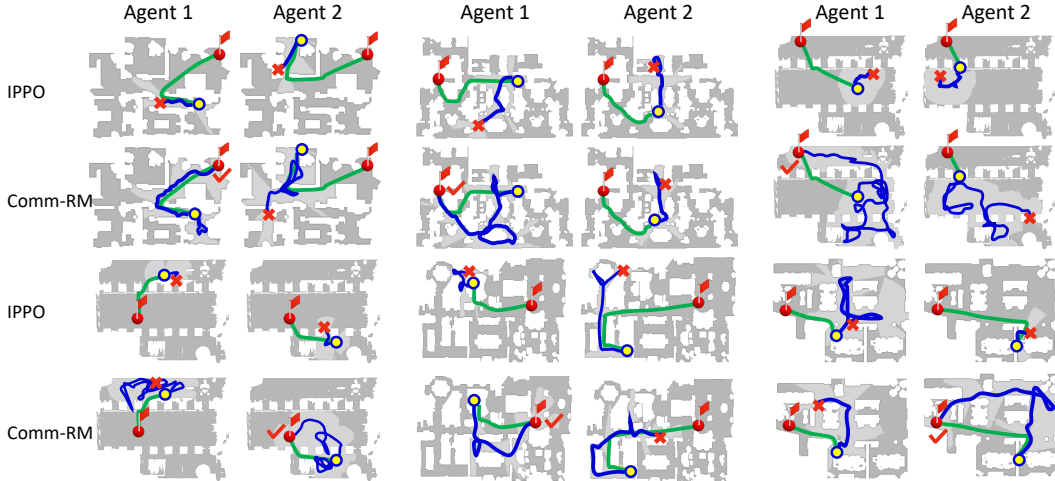


Figure 6: The trajectory visualization results of the IPPO agent and the communicative agent in the testing environment. The red circle with a red flag is the position where the target located. The yellow circle is the starting position of an agent. The green line indicates the shortest path, and the blue line is the actual navigation path. The red cross indicate the location where an agent fails.

291 **The Difficulty of Two Sub-tasks** The Fig. 5(b), (c) ablate the difficulty of two sub-tasks. We train  
 292 the IPPO baselines on individual-target task and shared-target task respectively. We find that the  
 293 individual-target task is significantly harder than the shared-target task, and the gap of difficulty is  
 294 increasing with more agent amounts. This experimentation result also proves the dataset analysis  
 295 result in Sec 3.3.

296 **Ablation for Communication** We train the model with historical communication mechanism, the  
 297 model with value communication mechanism and the IPPO baseline on 2 agents, 4 agents and  
 298 6 agents scenario. The result if shown in the Fig. 5(d), (e), (f), where the model with historical  
 299 communication mechanism significantly outperform other two models. In addition, we find that the  
 300 value communication mechanism cause overfitting in the MAIN task.

301 A more detailed comparison is shown in Tab. 2. We test our baseline models and the Comm-RM  
 302 model in both 2-agents and 3-agent scenarios. We find that the third variant, whose structure is shown  
 303 in Fig. 4, performs the best and largely outperforms other methods. We conclude from this figure  
 304 that communication mechanism is quite important for the MAIN task. A proper communication  
 305 mechanism largely improves the performance while a bad design of cooperative mechanism may  
 306 introduce noise or cause overfitting. Moreover, we find that the results of the IPPO model and the  
 307 single-agent PPO model tested in the multi-agent environment still competitive.

308 **Visualization for Navigation Process** In Fig. 6, we visualize the navigation process of two models:  
 309 the IPPO baseline model and the Comm-RM model. In this figure, at least one agent from the  
 310 communicative model successfully reaches the target. We find that the agents with cooperative  
 311 communication is able to explore larger area and navigation for a longer trajectory. Similar result is  
 312 also observed in Tab. 2. We find that the agents with communication tend to explore different areas  
 313 in a room. It indicates that the agents is able to learn to navigate seperately and communicate the  
 314 exploration result, which largely improve the navigation efficiency.

## 315 6 Conclusion

316 In this paper, we propose a novel Multi-Agent Indoor Navigation (MAIN) benchmark to research  
 317 on multi-agent problem in a realistic environment. We collect a large-scale dataset for researching  
 318 on MAIN and analysis the advantage of our dataset. We benchmark multiple baseline models in  
 319 MAIN and find that traditional MARL methods cannot solve MAIN due to the unique challenges in  
 320 MAIN such as little observation overlap and high variance of the embodied image view. By doing  
 321 experimentation, We discover that the model with historical communication message significantly  
 322 helps multi-agent navigation in MAIN. In the future, we are going to research on MARL problems  
 323 based on MAIN and keep updating the dataset and the codebase of MAIN.

324 **References**

- 325 [1] M. Aiello, E. Cambiaso, R. Canonico, L. Maccari, M. Mellia, A. Pescapè, and I. Vaccari. Ippo: A  
326 privacy-aware architecture for decentralized data-sharing. *arXiv:2001.06420*, 2020.
- 327 [2] P. Anderson, A. X. Chang, D. S. Chaplot, A. Dosovitskiy, S. Gupta, V. Koltun, J. Kosecka, J. Malik,  
328 R. Mottaghi, M. Savva, and A. R. Zamir. On evaluation of embodied navigation agents. *arXiv:1807.06757*,  
329 2018.
- 330 [3] P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson, N. Sunderhauf, I. Reid, S. Gould, and A. van den  
331 Hengel. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real  
332 environments. In *CVPR*, 2018.
- 333 [4] B. Baker, I. Kanitscheider, T. Markov, Y. Wu, G. Powell, B. McGrew, and I. Mordatch. Emergent tool use  
334 from multi-agent autotutorials. In *ICLR*, 2020.
- 335 [5] N. Bard, J. N. Foerster, S. Chandar, N. Burch, M. Lanctot, H. F. Song, E. Parisotto, V. Dumoulin, S. Moitra,  
336 E. Hughes, I. Dunning, S. Mourad, H. Larochelle, M. G. Bellemare, and M. Bowling. The hanabi challenge:  
337 A new frontier for ai research. *Artificial Intelligence*, 2020.
- 338 [6] D. Batra, A. Gokaslan, A. Kembhavi, O. Maksymets, R. Mottaghi, M. Savva, A. Toshev, and E. Wijnmans.  
339 Objectnav revisited: On evaluation of embodied agents navigating to objects. *arXiv:2006.13171*, 2020.
- 340 [7] C. Berner, G. Brockman, B. Chan, V. Cheung, P. Debiak, C. Dennison, D. Farhi, Q. Fischer, S. Hashme,  
341 C. Hesse, R. Józefowicz, S. Gray, C. Olsson, J. Pachocki, M. Petrov, H. P. de Oliveira Pinto, J. Raiman,  
342 T. Salimans, J. Schlatter, J. Schneider, S. Sidor, I. Sutskever, J. Tang, F. Wolski, and S. Zhang. Dota 2 with  
343 large scale deep reinforcement learning. *arXiv:1912.06680*, 2019.
- 344 [8] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang.  
345 Matterport3d: Learning from rgb-d data in indoor environments. *arXiv:1709.06158*, 2017.
- 346 [9] D. S. Chaplot, D. Gandhi, S. Gupta, A. Gupta, and R. Salakhutdinov. Learning to explore using active  
347 neural slam. In *ICLR*, 2020.
- 348 [10] D. S. Chaplot, R. Salakhutdinov, A. Gupta, and S. Gupta. Neural topological slam for visual navigation. In  
349 *CVPR*, 2020.
- 350 [11] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning  
351 phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv:1406.1078*,  
352 2014.
- 353 [12] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation for games, robotics and machine  
354 learning. <http://pybullet.org>, 2016–2021.
- 355 [13] C. S. de Witt, T. Gupta, D. Makoviichuk, V. Makoviyshuk, P. H. S. Torr, M. Sun, and S. Whiteson. Is  
356 independent learning all you need in the starcraft multi-agent challenge? *CoRR*, 2020.
- 357 [14] M. Deitke, W. Han, A. Herrasti, A. Kembhavi, E. Kolve, R. Mottaghi, J. Salvador, D. Schwenk, E. Van-  
358 derBilt, M. Wallingford, et al. Robothor: An open simulation-to-real embodied ai platform. In *CVPR*,  
359 2020.
- 360 [15] J. Foerster, I. A. Assael, N. de Freitas, and S. Whiteson. Learning to communicate with deep multi-agent  
361 reinforcement learning. In *NeurIPS*, 2016.
- 362 [16] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson. Counterfactual multi-agent policy  
363 gradients. In *AAAI*, 2018.
- 364 [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- 365 [18] S. Hu, F. Zhu, X. Chang, and X. Liang. Updet: Universal multi-agent rl via policy decoupling with  
366 transformers. In *ICLR*, 2021.
- 367 [19] M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu. Rein-  
368 forcement learning with unsupervised auxiliary tasks. In *ICLR*, 2016.
- 369 [20] L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: a survey. *Journal of Artificial*  
370 *Intelligence Research*, 1996.
- 371 [21] E. Kolve, R. Mottaghi, D. Gordon, Y. Zhu, A. Gupta, and A. Farhadi. Ai2-thor: An interactive 3d  
372 environment for visual ai. *arXiv:1712.05474*, 2017.

- 373 [22] L. Kraemer and B. Banerjee. Multi-agent reinforcement learning as a rehearsal for decentralized planning.  
374 *Neurocomputing*, 2016.
- 375 [23] J. Z. Leibo, V. Zambaldi, M. Lanctot, J. Marecki, and T. Graepel. Multi-agent reinforcement learning in  
376 sequential social dilemmas. In *AAMAS*, 2017.
- 377 [24] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous  
378 control with deep reinforcement learning. In *ICLR*, 2016.
- 379 [25] M. L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *ICML*, 1994.
- 380 [26] S. Liu, G. Lever, J. Merel, S. Tunyasuvunakool, N. Heess, and T. Graepel. Emergent coordination through  
381 competition. In *ICLR*, 2019.
- 382 [27] R. Lowe, Y. Wu, A. Tamar, J. Harb, O. P. Abbeel, and I. Mordatch. Multi-agent actor-critic for mixed  
383 cooperative-competitive environments. In *NeurIPS*, 2017.
- 384 [28] A. Mahajan, T. Rashid, M. Samvelyan, and S. Whiteson. Maven: Multi-agent variational exploration.  
385 *arXiv:1910.07483*, 2019.
- 386 [29] L. Matignon, L. Jeanpierre, and A.-I. Mouaddib. Coordinated multi-robot exploration under communication  
387 constraints using decentralized markov decision processes. In *AAAI*, 2012.
- 388 [30] P. Mirowski, R. Pascanu, F. Viola, H. Soyer, A. J. Ballard, A. Banino, M. Denil, R. Goroshin, L. Sifre,  
389 K. Kavukcuoglu, D. Kumaran, and R. Hadsell. Learning to navigate in complex environments. In *ICLR*,  
390 2016.
- 391 [31] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Harley, T. P. Lillicrap, D. Silver, and K. Kavukcuoglu.  
392 Asynchronous methods for deep reinforcement learning. In *ICML*, 2016.
- 393 [32] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. A. Riedmiller. Playing  
394 atari with deep reinforcement learning. *arXiv:1312.5602*, 2013.
- 395 [33] I. Mordatch and P. Abbeel. Emergence of grounded compositional language in multi-agent populations. In  
396 *AAAI*, 2017.
- 397 [34] A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transac-*  
398 *tions on Automatic Control*, 2009.
- 399 [35] F. A. Oliehoek and C. Amato. *A concise introduction to decentralized POMDPs*. Springer, 2016.
- 400 [36] F. A. Oliehoek, M. T. Spaan, and N. Vlassis. Optimal and approximate q-value functions for decentralized  
401 pomdps. *JAIR*, 2008.
- 402 [37] P. Paquette, Y. Lu, S. Bocco, M. O. Smith, S. Ortiz-Gagné, J. K. Kummerfeld, S. Singh, J. Pineau, and  
403 A. Courville. No press diplomacy: Modeling multi-agent gameplay. *arXiv:1909.02128*, 2019.
- 404 [38] P. Peng, Q. Yuan, Y. Wen, Y. Yang, Z. Tang, H. Long, and J. Wang. Multiagent bidirectionally-coordinated  
405 nets for learning to play starcraft combat games. *arXiv:1703.10069*, 2017.
- 406 [39] D. Pérez-Liébana, K. Hofmann, S. P. Mohanty, N. Kuno, A. Kramer, S. Devlin, R. D. Gaina, and D. Ionita.  
407 The multi-agent reinforcement learning in malmÖ (marlÖ) competition. *arXiv:1901.08129*, 2019.
- 408 [40] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson. Qmix: Monotonic  
409 value function factorisation for deep multi-agent reinforcement learning. In *ICML*, 2018.
- 410 [41] N. Sambasivan, S. Kapania, H. Highfill, D. Akrong, P. Paritosh, and L. M. Aroyo. “everyone wants to do  
411 the model work, not the data work”: Data cascades in high-stakes ai. In *CHI*, 2021.
- 412 [42] M. Samvelyan, T. Rashid, C. S. De Witt, G. Farquhar, N. Nardelli, T. G. Rudner, C.-M. Hung, P. H. Torr,  
413 J. Foerster, and S. Whiteson. The starcraft multi-agent challenge. *arXiv:1902.04043*, 2019.
- 414 [43] M. Savva, A. X. Chang, A. Dosovitskiy, T. A. Funkhouser, and V. Koltun. Minos: Multimodal indoor  
415 simulator for navigation in complex environments. *arXiv:1712.03931*, 2017.
- 416 [44] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik,  
417 et al. Habitat: A platform for embodied ai research. In *ICCV*, 2019.

- 418 [45] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using  
419 generalized advantage estimation. In *ICLR 2016 : International Conference on Learning Representations*  
420 *2016*, 2016.
- 421 [46] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms.  
422 *arXiv:1707.06347*, 2017.
- 423 [47] K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi. Qtran: Learning to factorize with transformation  
424 for cooperative multi-agent reinforcement learning. In *ICML*, 2019.
- 425 [48] S. Sukhbaatar, A. Szlam, and R. Fergus. Learning multiagent communication with backpropagation. In  
426 *NeurIPS*, 2016.
- 427 [49] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. Zambaldi, M. Jaderberg, M. Lanctot, N. Son-  
428 nerat, J. Z. Leibo, K. Tuyls, et al. Value-decomposition networks for cooperative multi-agent learning.  
429 *arXiv:1706.05296*, 2017.
- 430 [50] R. S. Sutton, D. A. McAllester, S. P. Singh, Y. Mansour, et al. Policy gradient methods for reinforcement  
431 learning with function approximation. In *NeurIPS*, 1999.
- 432 [51] M. Tan. Multi-agent reinforcement learning: independent vs. cooperative agents. In *ICML*, 1997.
- 433 [52] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell,  
434 T. Ewalds, P. Georgiev, J. Oh, D. Horgan, M. Kroiss, I. Danihelka, A. Huang, L. Sifre, T. Cai, J. P. Agapiou,  
435 M. Jaderberg, A. S. Vezhnevets, R. Leblond, T. Pohlen, V. Dalibard, D. Budden, Y. Sulsky, J. Molloy,  
436 T. L. Paine, Çağlar Gülçehre, Z. Wang, T. Pfaff, Y. Wu, R. Ring, D. Yogatama, D. Wünsch, K. McKinney,  
437 O. Smith, T. Schaul, T. P. Lillicrap, K. Kavukcuoglu, D. Hassabis, C. Apps, and D. Silver. Grandmaster  
438 level in starcraft ii using multi-agent reinforcement learning. *Nature*, 2019.
- 439 [53] O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, M. Yeo, A. Makhzani, H. Küttler, J. Aga-  
440 piou, J. Schrittwieser, et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv:1708.04782*,  
441 2017.
- 442 [54] R. E. Wang, M. Everett, and J. P. How. R-MADDPG for partially observable environments and limited  
443 communication. *CoRR*, 2020.
- 444 [55] T. Wang, T. Gupta, A. Mahajan, B. Peng, S. Whiteson, and C. Zhang. Rode: Learning roles to decompose  
445 multi-agent tasks. *arXiv:2010.01523*, 2020.
- 446 [56] X. Wang, Q. Huang, A. Celikyilmaz, J. Gao, D. Shen, Y.-F. Wang, W. Y. Wang, and L. Zhang. Reinforced  
447 cross-modal matching and self-supervised imitation learning for vision-language navigation. In *CVPR*,  
448 2018.
- 449 [57] S. Wani, S. Patel, U. Jain, A. X. Chang, and M. Savva. Multion: Benchmarking semantic map memory  
450 using multi-object navigation. *arXiv:2012.03912*, 2020.
- 451 [58] Y. Wu, Y. Wu, G. Gkioxari, and Y. Tian. Building generalizable agents with a realistic and rich 3d  
452 environment. *arXiv:1801.02209*, 2018.
- 453 [59] F. Xia, W. B. Shen, C. Li, P. Kasimbeg, M. E. Tchapmi, A. Toshev, R. Martin-Martin, and S. Savarese.  
454 Interactive gibbon benchmark: A benchmark for interactive navigation in cluttered environments. *IEEE*  
455 *Robotics and Automation Letters*, 5(2), 2020.
- 456 [60] F. Xia, A. R. Zamir, Z. He, A. Sax, J. Malik, and S. Savarese. Gibson env: Real-world perception for  
457 embodied agents. In *CVPR*, 2018.
- 458 [61] Y. Yang, Y. Wen, J. Wang, L. Chen, K. Shao, D. Mguni, and W. Zhang. Multi-agent determinantal  
459 q-learning. In *ICML*, 2020.
- 460 [62] C. Yu, A. Velu, E. Vinitsky, Y. Wang, A. M. Bayen, and Y. Wu. The surprising effectiveness of mappo in  
461 cooperative, multi-agent games. *arXiv:2103.01955*, 2021.
- 462 [63] C. Yu, A. Velu, E. Vinitsky, Y. Wang, A. M. Bayen, and Y. Wu. The surprising effectiveness of MAPPO in  
463 cooperative, multi-agent games. *CoRR*, 2021.
- 464 [64] J. Zhang, L. Tai, M. Liu, J. Boedecker, and W. Burgard. Neural slam: Learning to explore with external  
465 memory. *arXiv:1706.09520*, 2017.
- 466 [65] K. Zhang, Z. Yang, and T. Basar. Multi-agent reinforcement learning: A selective overview of theories and  
467 algorithms. *arXiv:1911.10635*, 2019.
- 468 [66] F. Zhu, Y. Zhu, X. Chang, and X. Liang. Vision-language navigation with self-supervised auxiliary  
469 reasoning tasks. In *CVPR*, 2020.

## 470 7 Checklist

- 471 1. For all authors...
- 472 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
473 contributions and scope? [Yes] See Section *Abstract* and *Introduction*.
- 474 (b) Have you read the ethics review guidelines and ensured that your paper conforms to  
475 them? [Yes]. We have read the ethics review guidelines and our paper is conforms to  
476 them.
- 477 (c) Did you discuss any potential negative societal impacts of your work?[Yes] Our paper  
478 is only for academic research purposes.
- 479 (d) Did you describe the limitations of your work? [Yes] Our experimentation is computa-  
480 tion costly, as shown in the implementation detail section.
- 481 2. If you are including theoretical results...
- 482 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 483 (b) Did you include complete proofs of all theoretical results? [N/A]
- 484 3. If you ran experiments...
- 485 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
486 mental results (either in the supplemental material or as a URL)? [Yes]
- 487 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were  
488 chosen)? [Yes]. We include the URL of our website in the abstract and supplementary  
489 materials.
- 490 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
491 ments multiple times)? [No] Due to the limit of computation resource, we have no  
492 time to investigate the error bars. We will make it up in the revision.
- 493 (d) Did you include the amount of compute and the type of resources used (e.g., type of  
494 GPUs, internal cluster, or cloud provider)? [Yes]
- 495 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 496 (a) If your work uses existing assets, did you cite the creators? [Yes] See Part *Experiment*.
- 497 (b) Did you mention the license of the assets? [Yes]. We use CC BY license that allows  
498 reusers to distribute, remix, adapt, and build upon the material in any medium or format,  
499 so long as attribution is given to the creator. The license allows for commercial use.
- 500 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes].  
501 We include a URL in the abstract that links to our website. Our website provide code  
502 and data to reproduce the results.
- 503 (d) Did you discuss whether and how consent was obtained from people whose data you're  
504 using/curating? [Yes] The data and methods are publicly released.
- 505 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
506 information or offensive content? [Yes] Data used in our work does not contain  
507 personally identifiable information or offensive content.
- 508 5. If you used crowdsourcing or conducted research with human subjects...
- 509 (a) Did you include the full text of instructions given to participants and screenshots, if  
510 applicable? [N/A]
- 511 (b) Did you describe any potential participant risks, with links to Institutional Review  
512 Board (IRB) approvals, if applicable? [N/A]
- 513 (c) Did you include the estimated hourly wage paid to participants and the total amount  
514 spent on participant compensation? [N/A]