

MA²E: ADDRESSING PARTIAL OBSERVABILITY IN MULTI-AGENT REINFORCEMENT LEARNING WITH MASKED AUTOENCODER

Anonymous authors

Paper under double-blind review

ABSTRACT

Centralized Training and Decentralized Execution (CTDE) is a widely adopted paradigm to solve cooperative multi-agent reinforcement learning (MARL) problems. Despite the successes achieved with CTDE, partial observability still limits cooperation among agents. While previous studies have attempted to overcome this challenge through communication, direct information exchanges could be restricted and introduce additional constraints. Alternatively, if an agent can infer the global information solely from local observations, it can obtain a global view without the need for communication. To this end, we propose the Multi-Agent Masked Auto-Encoder (MA²E), which utilizes the masked auto-encoder architecture to infer the information of other agents from partial observations. By employing masking to learn to reconstruct global information, MA²E serves as an inference module for individual agents within the CTDE framework. MA²E can be easily integrated into existing MARL algorithms and has been experimentally proven to be effective across a wide range of environments and algorithms.

1 INTRODUCTION

In cooperative multi-agent tasks, the environments are typically partially observable, where individual agents do not have complete access to the global information. For instance, consider a motivating scenario where allied and enemy forces are engaged in combat as illustrated in Figure 1. An ally agent 1 might perceive the current situation as favorable by observing only one enemy, while ally agents 2 and 3 observe four and five enemies, respectively, and therefore assess the situation as unfavorable. This discrepancy makes it difficult for the agents to make cooperative decisions. Therefore, addressing partial observability is a crucial challenge in the field of cooperative Multi-Agent Reinforcement Learning (MARL) (Nguyen et al., 2020; Canese et al., 2021; Ning & Xie, 2024).

One straightforward approach is fully centralized training and control of a joint policy for all agents which simplifies the MARL problem into a single-agent one. However, the centralized setting suffers from scalability and heavy computational costs as the number of agents increases due to the curse of dimensionality (Gronauer & Diepold, 2022; Du & Ding, 2021; Canese et al., 2021).

Centralized Training and Decentralized Execution (CTDE) paradigm provides a structured framework for mitigating aforementioned issues (Canese et al., 2021; Ning & Xie, 2024; Du & Ding, 2021; Gronauer & Diepold, 2022), but still fails to fully resolve the problem of partial observability. Although the global information is used during training, each agent relies solely on its local observations during execution. Such discrepancies between learning and execution can limit collaborative decision-making (Shao et al., 2022; Yuan et al., 2022; Guan et al., 2022).

Previous works have explored the strategies to relieve drawback of decentralized execution. Relaxing CTDE with communication by allowing direct messaging among agents have been extensively studied, but it faces limitations in real-world scenarios where inter-agent communication is not permitted or is constrained by various factors such as limited bandwidth and noisy channels (Zhu et al., 2022; Ning & Xie, 2024). Without communication, shared structures such as common knowledge (Schroeder de Witt et al., 2019) or an one-hot consensus (Xu et al., 2023) have been proposed to foster cooperation. However, such abstractly encoded representations may be lossy in complex environments.

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Masked modeling, which learns to reconstruct original data from masked input, has shown notable achievements in the language (Devlin et al., 2018; Liu, 2019; Clark, 2020), the vision (Chen et al., 2020; Bao et al., 2021; He et al., 2022), and even in the single-agent reinforcement learning (RL) domains (Seo et al., 2023; Liu et al., 2023a; Wu et al., 2023). Inspired by those successes, in this paper, we propose Multi-Agent Masked AutoEncoder (MA²E) which enables agents to infer the global information based solely on partial observations. Our key idea is to use a masked autoencoder (MAE) (He et al., 2022), which masks a subset of the input and is trained to recover the missing (masked) data using the visible (unmasked) regions, to reconstruct full trajectories of all agents from local observations. While such masked modeling usually targets to solve the downstream tasks with complete data after training from incomplete input, our focus is to utilize MAE’s capability to handle partially observable data to infer the entire information. In MA²E, MAE performs masking on a per-agent basis and reconstructs the original global input, allowing the inference of global information from partial observations. This enhances the agent’s ability to deduce the situations of other agents and aids in making appropriate decisions. Furthermore, by separately training MA²E from agents’ policies, MA²E can be easily plugged into existing MARL algorithms as backbone networks and the learned MAE can be transferred to other backbone policies.

We experimentally evaluate our approach on the Starcraft Multi-agent Challenge (SMAC) (Samvelyan et al., 2019), SMACv2 (Ellis et al., 2023), and Google Research Football (GRF) (Kurach et al., 2020) environments. The experimental results consistently demonstrate that MA²E achieves faster convergence and higher sample efficiency compared to fine-tuned QMIX (Hu et al., 2021), which is the state-of-the-art MARL algorithm. Additionally, MA²E shows comparable or superior performance compared to the cases where full observations are provided or communication is employed, substantiating the ability of MA²E to effectively infer full observations from partial observations. Furthermore, employing MA²E achieves performance gains for various value-based and policy-based MARL methods as a backbone algorithm. Finally, to effectively integrate MA²E into the backbone MARL algorithms, we propose an appropriate configuration for applying MA²E across different hyperparameters and masking strategies in ablation studies.

2 RELATED WORK

Partial Observability in MARL: Partial observability is a fundamental challenge in cooperative MARL. As a naïve adoption of a fully centralized setting results in an intractable computational complexity (Canese et al., 2021; Du & Ding, 2021; Gronauer & Diepold, 2022), alternative strategies have been explored. Centralized Training and Decentralized Execution (CTDE) is a popular framework to solve partial observability (Canese et al., 2021; Ning & Xie, 2024; Du & Ding, 2021; Gronauer & Diepold, 2022). In value-based CTDE algorithms, value decomposition which factorizes a joint value function into individual ones is dominant (Sunehag et al., 2017; Rashid et al., 2018; Lowe et al., 2017; Foerster et al., 2018; Yu et al., 2022). Another major branch is policy-based CTDE methods, where actor-critic approach is widely-adopted (Lowe et al., 2017; Foerster et al., 2018; Yu et al., 2022). We aim to alleviate the limitation of fully decentralized execution solely with local observations while following the CTDE paradigm. Subsequent to some pioneer works (Sukhbaatar et al., 2016; Foerster et al., 2016), communication among agents has been extensively studied in cooperative MARL (Yuan et al., 2022; Hu et al., 2021; 2024). Although communication can mitigate partial observability, it may not be possible in target environments or constrained by diverse factors (Zhu et al., 2022; Ning & Xie, 2024). Li et al. (2023) gradually shifts from communication to tacit collaboration, but it requires communication during exploration. Our goal is to remove the dependence on communication. Another line of research is to utilize shared structures dependent only on local information during

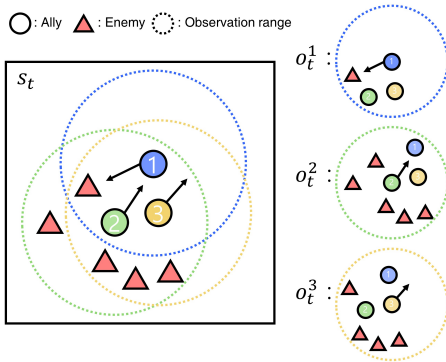


Figure 1: A motivating combat scenario. The solid circles and red triangles denote the ally agents and enemies, respectively. The dotted circles are observation ranges of each ally agent. While the current global state s_t is identical, each agent perceives the situation differently due to the discrepancy among their partial observations o_t^1, o_t^2, o_t^3 .

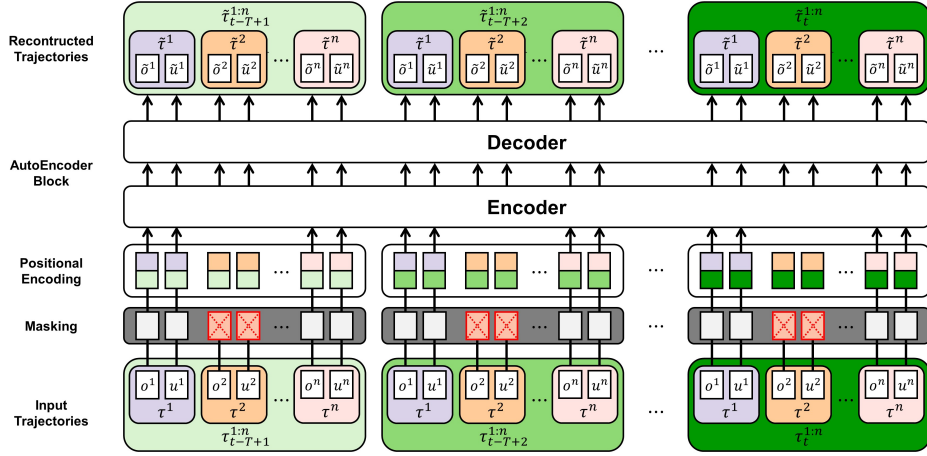


Figure 2: The architecture of MA²E. During centralized training, MA²E masks out k agents’ trajectories from all agents’ trajectories $\tau_{t-T+1:t}^n$ and learns to reconstruct them. The positional encoding is applied considering both the time and the agent information. During decentralized execution, trajectories of other agents are masked except the trajectories of current agent.

execution. MACKRL (Schroeder de Witt et al., 2019) employs common knowledge among groups of agents with a hierarchical policy tree. COLA (Xu et al., 2023) infers the same one-hot consensus for all agents from different local observations. Compared to those limited abstractions, our method enables each agent to recover the global trajectories of all agents from its local observations.

Masked Modeling in Reinforcement Learning: Masked modeling to reconstruct original data from masked ones emerges as a powerful technique in diverse domains including vision (Chen et al., 2020; Bao et al., 2021; He et al., 2022) and NLP (Devlin et al., 2018; Liu, 2019; Clark, 2020). Recently, the idea has been extended to the RL field. Masked World Model (Seo et al., 2023) is trained to reconstruct pixels from masked convolutional features to learn a latent dynamics model. MaskDP (Liu et al., 2023a) and MTM (Wu et al., 2023) apply masking to a portion of input trajectories and learn to reconstruct them. They show generalization abilities on diverse tasks by manipulating masking patterns desired for the target task. Aforementioned methods aim to solve the single-agent RL tasks. In contrast, we explore MARL problems with masked modeling. MA2CL (Song et al., 2023) deploys masked modeling with contrastive loss. The policy and the reconstruction module are jointly trained to obtain encoder representation for the policy with enhanced collaboration. In contrast, MA²E separately trains the policy and the reconstruction module, and its main focus is to recover global information rather than encoder representation. MaskMA (Liu et al., 2023b) treats MARL as a sequence modeling problem and masking is used to predict next action, rather than reconstructing the inputs. In contrast, our aim is to recover global information using MAE for decentralized execution.

3 METHOD

Problem Formulation: We consider a cooperative MARL task which can be described as a Dec-POMDP (Oliehoek et al., 2016). Dec-POMDP is defined as a tuple $G = \langle S, U, P, r, Z, O, n, \gamma \rangle$ where $s \in S$ denotes the global environment state. At each time step, each agent $i \in \{1, \dots, n\}$ chooses an action $u^i \in U$ which consists a joint action $\mathbf{u} = \{u^1, \dots, u^n\} \in U^n$. The environment follows the transition function $P(s'|s, \mathbf{u}) : S \times U^n \times S \rightarrow [0, 1]$. All the agents share the common reward function $r(s, \mathbf{u}) : S \times U^n \rightarrow \mathbb{R}$. Each agent only obtains an individual partial observation $o^i \in Z$ with the observation function $O(s, i) : S \times N \rightarrow Z$, and has an observation-action history

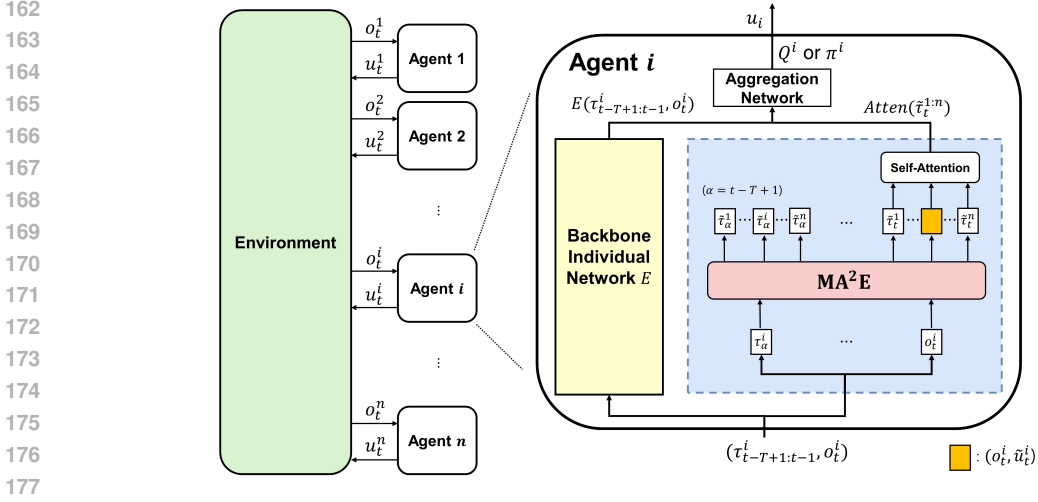


Figure 3: Incorporating MA²E into individual agents. Given a backbone MARL algorithm, MA²E is integrated into the backbone individual network of each agent during decentralized execution. MA²E serves as a global information inference module by only using local observations of the current agent and inferring other agents’ information.

$\tau^i \in T \equiv (Z \times U)^*$ and an individual policy $\pi^i(u^i | \tau^i) : T \times U \rightarrow [0, 1]$. The objective is to maximize the expected return $\mathbb{E}_{s_{t+1:\infty}, \mathbf{u}_{t+1:\infty}} [\sum_{m=0}^{\infty} \gamma^m r_{t+m} | s_t, \mathbf{u}_t]$ with a discount factor $\gamma \in [0, 1)$.

Masking: While random masking to states and actions is commonly adopted in single agent RL (Liu et al., 2023a; Wu et al., 2023), our aim is to infer other agents’ information only with local observations hence we apply masking at the agent level. During centralized training, we can access all agents’ trajectories $\tau_{t-T+1:t}^n = \{\tau_{t-T+1}^n, \dots, \tau_t^n\}$ where $\tau_t^i = (o_t^i, u_t^i)$. As illustrated in Figure 2, we randomly select k agents and mask out all trajectories of these agents. Then MA²E learns to recover the entire trajectories from the masked trajectories that it outputs $\tilde{\tau}_{t-T+1:t}^{1:n}$ where $\tilde{\tau}_t^i$ represents a recovered trajectory for the agent i at time step t .

In decentralized execution (for exploration during training or deployment after training), apart from the local observations and actions of the agent $(\tau_{t-T+1:t-1}^i, o_t^i)$, other information such as observations from other agents cannot be obtained. The agent only utilizes its local observations for the action selection and MA²E infers global information from local observations by restoring masked areas during execution, as depicted in Figure 3. The comparison between agent level and random masking strategies can be found in Section 4.4 and [Appendix H.4](#).

Architecture: The detailed architecture of MA²E is described in Figure 2. MA²E follows an auto-encoder architecture (Bank et al., 2020), with the encoder and decoder structured similarly to Transformer (Vaswani et al., 2017), using Multi-Head Attention (MHA) and feedforward networks. The decoder does not use masked MHA. Instead, a separate layer is utilized for masking before the input data is fed into the encoder and decoder. In other words, observations and actions are embedded in the embedding layer, undergo masking and positional encoding layers, then are fed into the encoder and decoder. For positional encoding, the orders of both time steps and agents are valuable information as MA²E takes histories of multiple agents. Therefore, we apply positional encoding considering *Agent* and *Time* dimensions, and the related results are presented in Section 4.4 and [Appendix H.4](#).

Incorporating MA²E into Individual Agent: MA²E is integrated into individual agents as illustrated in Figure 3. Following this procedure, we note that MA²E can be easily plugged into existing MARL algorithms. Depending on the selected MARL algorithm, each agent has a backbone network E corresponding to the value function Q^i or policy π^i . MA²E is integrated internally within each agent in addition to E . Each agent only takes its own local observations to determine actions; hence, inputs for both E and MA²E are also local observations. Then MA²E reconstructs trajectories of all agents from its local trajectories, where the output for the agent i at the current time t is replaced by (o_t^i, \tilde{u}_t^i) since o_t^i is already observed. The recovered trajectory for the current time t , $\tilde{\tau}_t^{1:n} = (\tilde{\tau}_t^1, \dots, \tilde{\tau}_t^n)$, is

fed into the self-attention layers. The self-attention layers use a self-attention mechanism (Vaswani et al., 2017) to selectively focus on other agents’ information for the agent’s current situation $\tilde{\tau}_t^i$. For the agent i and j , the attention weight $w_{i,j}$ is calculated as:

$$w_{i,j} = \frac{\exp(\tilde{\tau}_t^i W_k^T W_q^T \tilde{\tau}_t^j)}{\sum_{j=1, \dots, i-1}^{i+1, \dots, n} \exp(\tilde{\tau}_t^i W_k^T W_q^T \tilde{\tau}_t^j)} \quad (1)$$

W_q and W_k are weight matrices such as fully connected layers for the query q and key k , respectively. Then a relevance-weighted value is computed by taking the weighted sum of the attention weights and values as shown below:

$$Atten_i(\tilde{\tau}_t^{1:n}) = \sum_{j=1, \dots, i-1}^{i+1, \dots, n} w_{i,j} W_v \tilde{\tau}_t^j, \quad (2)$$

where W_v is another weight matrix. From the backbone network E , we take the outputs before the last layer $E(\tau_{t-T+1:t-1}^i, o_t^i)$ as latent variables, and both $E(\tau_{t-T+1:t-1}^i, o_t^i)$ and $Atten_i(\tilde{\tau}_t^{1:n})$ are fed into the aggregation network. Finally, the aggregation network produces the individual value Q^i or policy π^i from which the agent chooses the action u^i . This framework enables fully decentralized execution for each agent without access to global information.

Training: The training of MA²E is divided into two stages: MA²E pre-training and MA²E fine-tuning. During the pre-training stage, MA²E is trained using samples collected by a random policy before the policy training begins. If we train MA²E from scratch and agents’ policy concurrently, the generated information from MA²E would be a noise that interferes with the policy training. Hence, the pre-training is needed to ensure that MA²E can properly infer unobserved information. The objective of MA²E, defined in Eq. (3), is to minimize the Mean Squared Error (MSE) loss between masked and true histories. Pre-training continues until the loss falls below a specified threshold.

$$L_{MA^2E} = \frac{1}{nT} \sum_{t=1}^T \sum_{i=1}^n (\tau_t^i - \tilde{\tau}_t^i)^2. \quad (3)$$

After pre-training of MA²E, the policy training begins. In this stage, fine-tuning of MA²E also takes place. However, the policy is independently trained according to the chosen MARL method, and MA²E is periodically updated with the loss defined in Eq. (3) after updating the policy sufficient times. Both the agents’ policy and MA²E are trained using collected trajectories in a replay buffer.

4 EXPERIMENT AND RESULT

Environment and Setup We evaluate MA²E on the StarCraft Multi-agent Challenge (SMAC) (Samvelyan et al., 2019), SMACv2 (Ellis et al., 2023), and Google Research Football (GRF) (Kurach et al., 2020) environments. SMAC is one of the most popular MARL benchmark, which covers a wide range of cooperative microcontrol scenarios. We conduct the experiments on SMAC HARD and SuperHARD scenarios. SMACv2 complements the deterministic property of SMAC by randomizing start positions and unit types, and changing the units’ sight and attack ranges. In GRF, multiple agents cooperate to play a football game. All the experiments are conducted during 2×10^6 time steps for each run, and we report the average win rates with the shaded standard error from three different random seeds. More experimental details can be found in Appendix A.

4.1 IMPROVING PERFORMANCE BY INCORPORATING MA²E INTO MARL ALGORITHMS

Firstly, we investigate the effectiveness of integrating MA²E into existing MARL methods to improve the performance. We employ a fine-tuned QMIX (Hu et al., 2021), which achieves SOTA performance through parameter fine-tuning from the vanilla QMIX (Rashid et al., 2018), as a backbone algorithm to apply MA²E. We compare MA²E with various MARL algorithms including VDN (Sunehag et al., 2017), QMIX (Rashid et al., 2018), OW-QMIX (Rashid et al., 2020), CW-QMIX (Rashid et al., 2020), QPLEX (Wang et al., 2020), MAC2L (Song et al., 2023), COLA (Xu et al., 2023) as baselines.

The performance comparison between the model with MA²E and the baselines is illustrated in Figure 4 and Figure 5. **QMIX+MA²E** is our proposed model, which is a fine-tuned QMIX (Hu

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

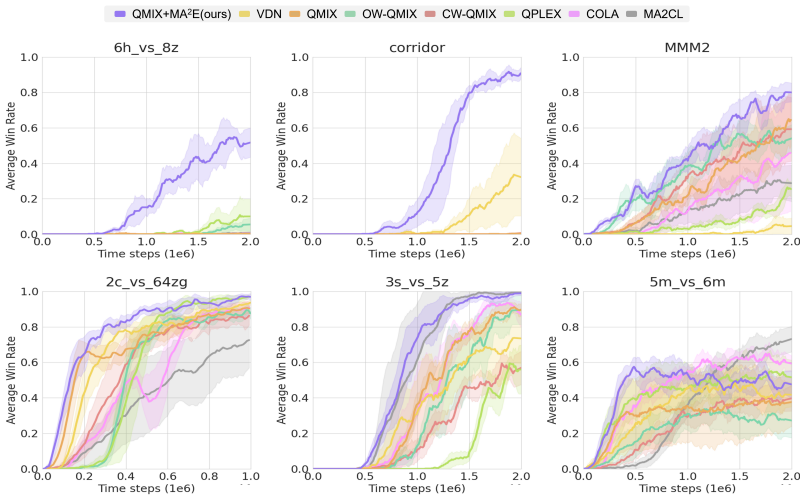


Figure 4: Performance comparison with baselines in SMAC scenarios. The blue line represents the model where MA²E is plugged into QMIX, while the other colored lines correspond to baselines. The proposed model performs better in both HARD and SuperHARD scenarios in SMAC.

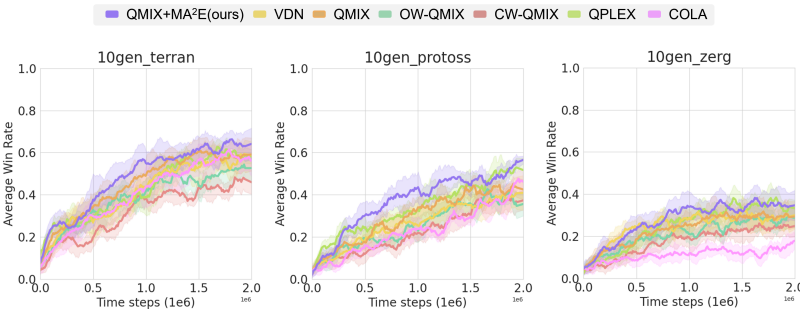


Figure 5: Performance comparison with baselines in SMACv2 scenarios. The blue line represents the model where MA²E is plugged into QMIX, while the other colored lines correspond to baselines.

et al., 2021) with the addition of MA²E. In the HARD and SuperHARD scenarios of SMAC and in SMACv2, QMIX+MA²E exhibits higher sample efficiency and superior win rates compare to all baselines. In scenarios with higher difficulty levels, such as *corridor* and *6h_vs_8z*, the performance difference between the proposed model and the baselines becomes more pronounced compared to relatively easier scenarios like *2c_vs_64zg*. Specifically, in challenging scenarios such as *corridor* and *6h_vs_8z*, baselines struggle to achieve victories even after two million time steps, whereas QMIX+MA²E attains high win rates. In addition, even when compared to QMIX, which is a fine-tuned version showing state-of-the-art performance in MARL, the proposed method shows superior performance. Even in the *10gen_terrain* scenario in SMACv2, QMIX+MA²E demonstrates higher performance compared to the baselines. [Additional results are in Appendix H.](#)

4.2 COMPARISON WITH THE FULL STATE AND COMMUNICATION METHODS

We evaluate whether MA²E can accurately infer full observations from partial observations. We compare the results of the model when using full state with the results of the model using MA²E. The fully observable model replaces the MA²E of individual agents which is Figure 3 with full observation and utilizes the output obtained through the self-attention layer. The results are depicted in Figure 6. The comparison across four SMAC scenarios shows a similar pattern between using full observation and using MA²E. Moreover, it is evident that the performance is better when using full observation or using MA²E compared to using only partial observations. Especially in scenarios like

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

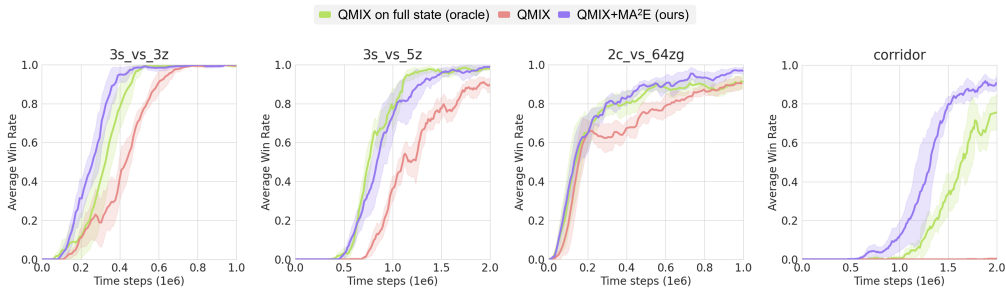


Figure 6: Performance comparison between using full observation and utilizing MA²E in different SMAC scenarios. Green line represents the case where QMIX is used with the full state, the blue line corresponds to our proposed model, and the red line represents the fine-tuned QMIX which is a baseline model in the experiment. The win rates of the model using the full state and the model with MA²E applied exhibit similar patterns.

Table 1: Comparison of MA²E observation inference accuracy based on training progress

Type	Before Pre-Training MA ² E	After Pre-Training MA ² E	After Fine-Tuning MA ² E
<ul style="list-style-type: none"> ● : True Allies ▲ : True Enemies ● : Inferred Allies ▲ : Inferred Enemies 			
Discrepancy			
Distance (L2 norm)	0.9933	0.3989	0.3588
Health (L1 norm)	0.6606	0.5585	0.4766
All Info (L1 norm)	0.5661	0.3984	0.3584

3s_vs_5z, where inter-agent information sharing is crucial and the difficulty is higher compared to other scenarios, the differences become more pronounced.

Table 1 compares the inference accuracy according to the training progress of MA²E. It shows the results of agent 1 inferring the observations of agent 2 in the 3s_vs_5z scenario in SMAC. From left to right in the table, the results are from models at progressively advanced stages of training, and the figure displays the inferred relative positions of agents from the observation inference results of other agents, plotted on a two-dimensional plane. In the figure, red and blue icons represent true values and the inferred values, respectively, indicating that as learning progresses, the inferred values become closer to the actual values. In addition, the values in the table mean the differences between the actual values and the MA²E’s deduced values for position, health, and the entire observation values of other agents. The discrepancies decrease as MA²E learns, demonstrating that MA²E can infer values closer to the real values through training. Based on the results, we can conclude that MA²E can successfully infer information similar to full states.

Another way to utilize information over partial observations is information exchange among agents through communication. We compare MA²E with communication methods in MARL domain including MAIC (Yuan et al., 2022), QMIX-att (Hu et al., 2021) and CommFormer (Hu et al., 2024). [The learning curves can be found in Appendix H.3](#). As shown in Table 2, QMIX+MA²E demonstrates performance that is comparable to or better than communication methods. Since MA²E does not rely on direct information exchange, its performance highlights the strengths of our method.

4.3 EXTENSIBILITY AND TRANSFERABILITY OF MA²E

To explore the applicability of MA²E to various MARL algorithms in a general context, we conduct experiments by applying MA²E to various value-based models such as VDN and QPLEX. The experimental results are as shown in the Figure 7. The results across scenarios in both SMAC and

Table 2: Comparison of MA²E with communication-based methods

Scenario	Steps	QMIX+MA ² E (ours)	MAIC	QMIX-att	CommFormer
3s_vs_5z	2M	0.99 ± 0.01	0.86 ± 0.05	0.41 ± 0.37	0.06 ± 0.04
2c_vs_64zg	1M	0.97 ± 0.03	0.76 ± 0.18	0.34 ± 0.31	0.59 ± 0.06
5m_vs_6m	2M	0.48 ± 0.08	0.69 ± 0.07	0.81 ± 0.04	0.61 ± 0.05
corridor	2M	0.91 ± 0.10	0.0 ± 0.0	0.30 ± 0.23	0.0 ± 0.0
6h_vs_8z	2M	0.53 ± 0.08	0.0 ± 0.0	0.14 ± 0.04	0.0 ± 0.0
MMM2	2M	0.79 ± 0.08	0.97 ± 0.07	0.20 ± 0.08	0.0 ± 0.0

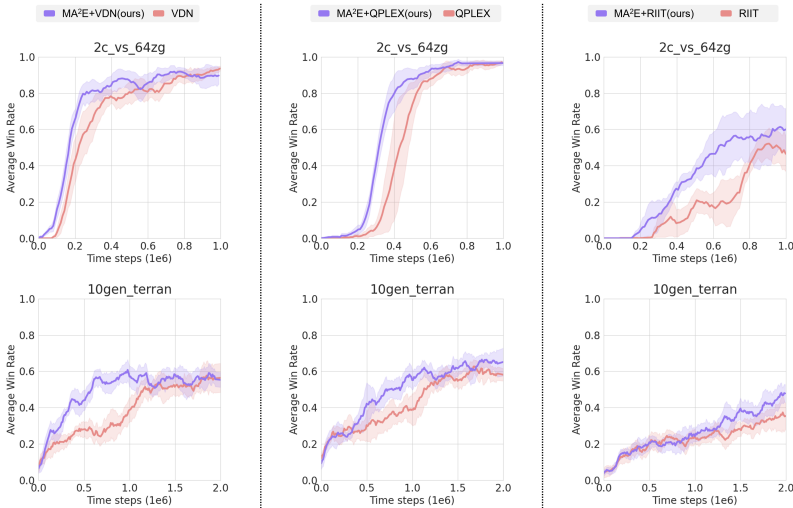


Figure 7: Performance comparison between the model incorporating MA²E in QPLEX and VDN and baselines in SMAC and SMACv2 scenarios. **(Left)**: Comparing the performance of the model with MA²E applied to VDN and the performance of VDN. **(Middle)**: Comparing the performance of the model with MA²E applied to QPLEX and the performance of QPLEX. **(Right)**: Comparing the performance of the model with MA²E applied to RIIT and the performance of RIIT to assess whether MA²E leads to performance improvement when applied to a policy-based method.

SMACv2 indicate that the addition of MA²E improves performance compared to the baselines. MA²E not only enhances sample efficiency but also converges to higher win rates in fewer time steps.

Furthermore, to confirm the applicability of MA²E to policy-based algorithms, we conduct experiments by incorporating MA²E into the policy-based algorithm RIIT (Hu et al., 2021). The right figures in Figure 7 illustrate that RIIT with MA²E exhibits better performance compared to the baseline models. The experimental results demonstrate that MA²E can be seamlessly integrated into both value-based and policy-based algorithms.

In order to confirm whether MA²E can be transferred to other backbone networks, MA²E is first trained with QMIX. Then the trained MA²E is transferred to the VDN network and only MA²E finetuning is performed as training progresses, without pre-training. Figure 8 compare the transferred MA²E with the MA²E that goes through the pre-training process. The model equipped with transferred MA²E performs better in the HARD scenario and shows comparable performance in the SuperHARD scenario, which demonstrates the transferability of MA²E.

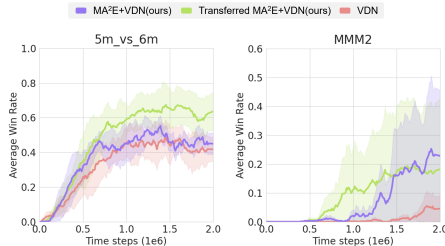


Figure 8: Performance comparison between transferred MA²E and default MA²E in different SMAC scenarios.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

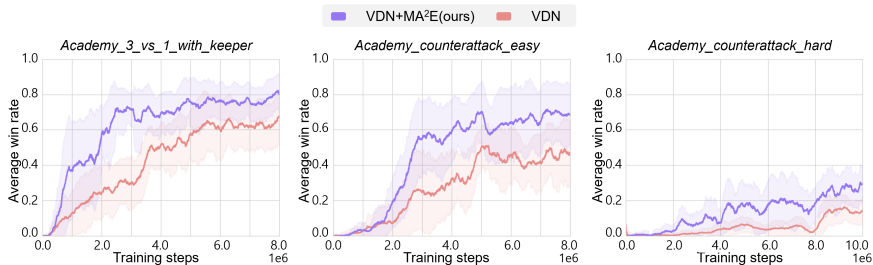


Figure 9: Performance comparison between the model incorporating MA²E in VDN and baselines in Google Research Football scenarios

To further verify the effectiveness of MA²E across diverse environments, additional experiments are conducted in GRF. We apply MA²E to VDN. As in Figure 9, integrating MA²E significantly improves performance in three different scenarios. Consequently, the above results demonstrate the extensibility of MA²E for various MARL algorithms and domains. We evaluated the performance of VDN, QPLEX, RIIT, and another policy-based algorithm, IPPO (de Witt et al., 2020) with MA²E in diverse scenarios and the results are illustrated in Appendix C. MA²E consistently improves the performance across different algorithms and tasks.

4.4 ABLATION STUDIES

To identify a suitable structure and hyper-parameters of MA²E, we conduct ablation studies across various settings. The comparison is conducted in the *3s_vs_5z* scenario of SMAC. The experimental results for comparison are illustrated in Figure 10 and in Appendix H.4.

Masking Strategy: We compare two masking strategies: agent-based and random masking. Random masking removes a random portion of the input data, irrespective of the agents. The masking ratio is randomly chosen from 0.15, 0.35, 0.5, 0.75, and 0.95, with the masking point also being randomly determined. This ratio follows the suggestion from Liu et al. (2023a). Figure 10 (a) compares win rates based on different masking strategies and agent-based masking is superior because it can effectively capture correlations between agents.

The Number of Trajectories: To test the appropriate trajectory length, we conduct a comparison of the number of trajectories in five different settings : 1, 3, 5, 7, and 9. Figure 10 (b) compares performance based on the number of trajectories. The performance does not show a consistent difference among the different number of trajectories. However, it can be observed that the performance is at its best when the number of trajectories is 5, and it deteriorates as the number becomes smaller or larger than 5. Therefore, to achieve good performance with MA²E, it is necessary to set appropriate number of trajectories depending on the scenario or environment.

Positional Encoding: To test the effect of the positional encoding design, we compare positional encoding using three different encoding methods : *Agent*, *Time*, and *Both*. *Agent* encoding distinguishes encoding based on agents, while *Time* encoding sets different values based on the trajectories. In the *Both* setting, encoding takes both *Agent* and *Time* into account. For example, half of the embedding space is dedicated to agent-based encoding, and the remaining half is utilized for time-based encoding. Figure 10 (c) compares performance based on the positional encoding design. As we can see, the performance is significantly better when considering both time and agent, as opposed to either of them individually. When considering only one of either time or agent, there is not a significant difference in performance compared to the case where MA²E is not used. In particular, when considering only time, there are intervals where the performance is even worse than without MA²E.

MA²E Training: As mentioned in Section 3, MA²E training is divided into two stages. To verify the importance of both stages, we report the performance of MA²E without pre-training (*w/o PT*) and without fine-tuning (*w/o FT*). As shown in Figure 10 (d), employing both pre-training and fine-tuning achieves the best performance. Moreover, there is a tendency for the performance of the case without pre-training (*w/o PT*) to be better than the case with pre-training only (*w/o FT*), but the variance of *w/o PT* is very high. This indicates that when MA²E is underfitted and its outputs

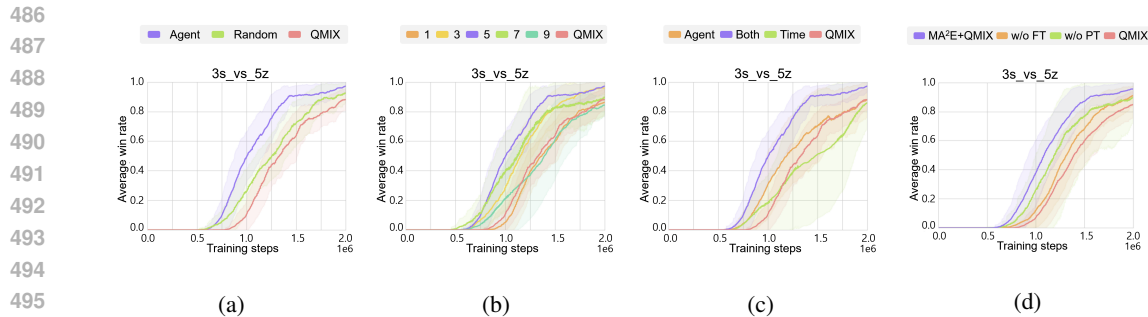


Figure 10: Performance comparison based on hyper parameter settings or strategies. **(a)**: Performance comparison based on masking strategies. The agent-based masking strategy applied model outperforms both the random masking strategy applied model and the default model. **(b)**: Performance comparison of MA²E based on the number of trajectories. **(c)**: Performance comparison of MA²E based on the positional encoding strategies. When considering both agent and time, it shows the best performance. **(d)**: Performance comparison of MA²E with or without pre-training and fine-tuning.

used in decision making are likely to act as noise. Therefore, pre-training is necessary to reliably improve the performance and fine-tuning is required to align MA²E with the policy.

5 CONCLUSION AND LIMITATIONS

In this paper, we propose a novel method MA²E, which utilizes the masked auto-encoder to address partial observability in the context of multi-agent reinforcement learning. MA²E empowers agents to derive global insights solely from own local information, thereby enhancing their collaborative decision-making capabilities. The proposed method seamlessly integrates with both value-based and policy-based MARL algorithms. Through extensive experimentation, we have substantiated its effectiveness in enhancing sample efficiency and elevating task-solving proficiency across a diverse set of scenarios within the SMAC, SMACv2, and GRF environments. We posit that the incorporation of MA²E into MARL, which extends the observational horizon of agents, stands as a pivotal advancement in MARL research. By introducing a methodology that addresses the challenge of partial observability in MARL, MA²E has emerged as a cornerstone in the field.

Nonetheless, our proposed model exhibits limitations. When agents are positioned far apart such that there is no overlap in their observations, the capacity to infer information about other agents is severely restricted, potentially resulting in the inference outcomes of MA²E being not useful. A more comprehensive delineation of these limitations can be found in Appendix E. Additionally, the structural characteristics of MA²E make it relatively difficult to scale to a varying number of agents. In future work, we plan to investigate approaches that enable scalability, such as team composition.

REPRODUCIBILITY STATEMENT

For the details of environments and hyperparameters, please refer Section 4 and Appendix A. To run our method, please download the supplementary material and follow the instructions in README files. We employed pymarl2 (Hu et al., 2021) or the official codes from the authors for baselines.

REFERENCES

- Dor Bank, Noam Koenigstein, and Raja Giryes. Autoencoders. *CoRR*, abs/2003.05991, 2020. URL <https://arxiv.org/abs/2003.05991>.
- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.
- Lorenzo Canese, Gian Carlo Cardarilli, Luca Di Nunzio, Rocco Fazzolari, Daniele Giardino, Marco Re, and Sergio Spanò. Multi-agent reinforcement learning: A review of challenges and applications.

- 540 *Applied Sciences*, 11(11), 2021. ISSN 2076-3417. doi: 10.3390/app11114948. URL <https://www.mdpi.com/2076-3417/11/11/4948>.
- 541
- 542
- 543 Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever.
- 544 Generative pretraining from pixels. In *International conference on machine learning*, pp. 1691–
- 545 1703. PMLR, 2020.
- 546
- 547 K Clark. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint*
- 548 *arXiv:2003.10555*, 2020.
- 549
- 550 Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS
- 551 Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft
- 552 multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- 553
- 554 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
- 555 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 556
- 557 Wei Du and Shifei Ding. A survey on multi-agent deep reinforcement learning: from the perspective
- 558 of challenges and applications. *Artificial Intelligence Review*, 54(5):3215–3238, 2021.
- 559
- 560 Benjamin Ellis, Jonathan Cook, Skander Moalla, Mikayel Samvelyan, Mingfei Sun, Anuj Ma-
- 561 hajan, Jakob Nicolaus Foerster, and Shimon Whiteson. SMACv2: An improved benchmark
- 562 for cooperative multi-agent reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=5OjLGiJW3u>.
- 563
- 564 Jakob Foerster, Ioannis Alexandros Assael, Nando De Freitas, and Shimon Whiteson. Learning
- 565 to communicate with deep multi-agent reinforcement learning. *Advances in neural information processing systems*, 29, 2016.
- 566
- 567 Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson.
- 568 Counterfactual multi-agent policy gradients. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- 569
- 570 Sven Gronauer and Klaus Diepold. Multi-agent deep reinforcement learning: a survey. *Artificial*
- 571 *Intelligence Review*, 55(2):895–943, 2022.
- 572
- 573 Cong Guan, Feng Chen, Lei Yuan, Chenghe Wang, Hao Yin, Zongzhang Zhang, and Yang Yu.
- 574 Efficient multi-agent communication via self-supervised information aggregation. *Advances in*
- 575 *Neural Information Processing Systems*, 35:1020–1033, 2022.
- 576
- 577 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
- 578 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- 579
- 580 Jian Hu, Haibin Wu, Seth Austin Harding, Siyang Jiang, and Shih-wei Liao. RIIT: rethinking the
- 581 importance of implementation tricks in multi-agent reinforcement learning. *CoRR*, abs/2102.03479,
- 582 2021. URL <https://arxiv.org/abs/2102.03479>.
- 583
- 584 Shengchao Hu, Li Shen, Ya Zhang, and Dacheng Tao. Learning multi-agent communication from
- 585 graph modeling perspective. *arXiv preprint arXiv:2405.08550*, 2024.
- 586
- 587 Karol Kurach, Anton Raichuk, Piotr Stańczyk, Michał Zajac, Olivier Bachem, Lasse Espeholt, Carlos
- 588 Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. Google research football:
- 589 A novel reinforcement learning environment. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 4501–4510, 2020.
- 590
- 591 Dapeng Li, Zhiwei Xu, Bin Zhang, and Guoliang Fan. From explicit communication to tacit
- 592 cooperation: A novel paradigm for cooperative marl. *arXiv preprint arXiv:2304.14656*, 2023.
- 593
- Fangchen Liu, Hao Liu, Aditya Grover, and Pieter Abbeel. Masked autoencoding for scalable and generalizable decision making. *Advances in neural information processing systems*, 2023a.

- 594 Jie Liu, Yinmin Zhang, Chuming Li, Chao Yang, Yaodong Yang, Yu Liu, and Wanli Ouyang. Masked
595 pretraining for multi-agent decision making, 2023b.
- 596
597 Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint*
598 *arXiv:1907.11692*, 2019.
- 599 Yong Liu, Weixun Wang, Yujing Hu, Jianye Hao, Xingguo Chen, and Yang Gao. Multi-agent game
600 abstraction via graph attention neural network. In *Proceedings of the AAAI conference on artificial*
601 *intelligence*, volume 34, pp. 7211–7218, 2020.
- 602
603 Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent
604 actor-critic for mixed cooperative-competitive environments. *Advances in neural information*
605 *processing systems*, 30, 2017.
- 606
607 Thanh Thi Nguyen, Ngoc Duy Nguyen, and Saeid Nahavandi. Deep reinforcement learning for
608 multiagent systems: A review of challenges, solutions, and applications. *IEEE Transactions on*
609 *Cybernetics*, 50(9):3826–3839, 2020. doi: 10.1109/TCYB.2020.2977374.
- 610
611 Zepeng Ning and Lihua Xie. A survey on multi-agent reinforcement learning and its ap-
612 plication. *Journal of Automation and Intelligence*, 3(2):73–91, 2024. ISSN 2949-8554.
613 doi: <https://doi.org/10.1016/j.jai.2024.02.003>. URL <https://www.sciencedirect.com/science/article/pii/S2949855424000042>.
- 614
615 Frans A Oliehoek, Christopher Amato, et al. *A concise introduction to decentralized POMDPs*,
616 volume 1. Springer, 2016.
- 617
618 Tabish Rashid, Mikayel Samvelyan, Christian Schröder de Witt, Gregory Farquhar, Jakob N. Foerster,
619 and Shimon Whiteson. QMIX: monotonic value function factorisation for deep multi-agent
620 reinforcement learning. *CoRR*, abs/1803.11485, 2018. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1803.11485)
621 [1803.11485](http://arxiv.org/abs/1803.11485).
- 622
623 Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted QMIX: expanding
624 monotonic value function factorisation. *CoRR*, abs/2006.10800, 2020. URL [https://arxiv.](https://arxiv.org/abs/2006.10800)
625 [org/abs/2006.10800](https://arxiv.org/abs/2006.10800).
- 626
627 Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli,
628 Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Foerster, and Shimon Whiteson. The
629 StarCraft Multi-Agent Challenge. *CoRR*, abs/1902.04043, 2019.
- 630
631 Christian Schroeder de Witt, Jakob Foerster, Gregory Farquhar, Philip Torr, Wendelin Boehmer, and
632 Shimon Whiteson. Multi-agent common knowledge reinforcement learning. *Advances in neural*
633 *information processing systems*, 32, 2019.
- 634
635 Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and Pieter Abbeel.
636 Masked world models for visual control. In *Conference on Robot Learning*, pp. 1332–1344. PMLR,
637 2023.
- 638
639 Jianzhun Shao, Zhiqiang Lou, Hongchang Zhang, Yuhang Jiang, Shuncheng He, and Xiangyang
640 Ji. Self-organized group for cooperative multi-agent reinforcement learning. *Advances in Neural*
641 *Information Processing Systems*, 35:5711–5723, 2022.
- 642
643 Haolin Song, Mingxiao Feng, Wengang Zhou, and Houqiang Li. Ma2cl: Masked attentive contrastive
644 learning for multi-agent reinforcement learning. *arXiv preprint arXiv:2306.02006*, 2023.
- 645
646 Sainbayar Sukhbaatar, Rob Fergus, et al. Learning multiagent communication with backpropagation.
647 *Advances in neural information processing systems*, 29, 2016.
- 648
649 Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi,
650 Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel.
651 Value-decomposition networks for cooperative multi-agent learning, 2017.
- 652
653 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
654 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
655 *systems*, 30, 2017.

648 Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. QPLEX: duplex dueling
649 multi-agent q-learning. *CoRR*, abs/2008.01062, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2008.01062)
650 [2008.01062](https://arxiv.org/abs/2008.01062).
651
652 Philipp Wu, Arjun Majumdar, Kevin Stone, Yixin Lin, Igor Mordatch, Pieter Abbeel, and Aravind
653 Rajeswaran. Masked trajectory models for prediction, representation, and control. *arXiv preprint*
654 *arXiv:2305.02968*, 2023.
655
656 Zhiwei Xu, Bin Zhang, Dapeng Li, Zeren Zhang, Guangchong Zhou, Hao Chen, and Guoliang Fan.
657 Consensus learning for cooperative multi-agent reinforcement learning. In *Proceedings of the*
658 *AAAI Conference on Artificial Intelligence*, volume 37, pp. 11726–11734, 2023.
659
660 Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The
661 surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information*
662 *Processing Systems*, 35:24611–24624, 2022.
663
664 Lei Yuan, Jianhao Wang, Fuxiang Zhang, Chenghe Wang, Zongzhang Zhang, Yang Yu, and Chongjie
665 Zhang. Multi-agent incentive communication via decentralized teammate modeling. In *Proceedings*
666 *of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 9466–9474, 2022.
667
668 Changxi Zhu, Mehdi Dastani, and Shihan Wang. A survey of multi-agent reinforcement learning
669 with communication. *arXiv preprint arXiv:2203.08975*, 2022.
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A EXPERIMENTAL DETAILS

In this section, we introduce the environments used in the experiments, the baseline algorithms, as well as the hyperparameters and computational resources. Experiments are carried out on NVIDIA A6000 and GTX3090 GPUs and AMD EPYC 7313 CPU.

A.1 ENVIRONMENTS

We conduct experiments in the following environments:

- StarCraft Multi-Agent Challenge (SMAC) (Samvelyan et al., 2019) from <https://github.com/oxwhirl/smac> which is licensed under MIT license.
- SMACv2 (Ellis et al., 2023) from <https://github.com/oxwhirl/smacv2> which is licensed under MIT license.
- Google Research Football (GRF) (Kurach et al., 2020) from <https://github.com/google-research/football> which is licensed under Apache License 2.0.

All algorithms are implemented based on the open-source framework *pymar12* (Hu et al., 2021) from <https://github.com/hijkzzz/pymar12> which is an augmented version of *pymar1* from <https://github.com/oxwhirl/pymar1>. Both are licensed under Apache License 2.0.

A.1.1 SMAC

The StarCraft Multi-Agent Challenge (SMAC) is one of the benchmarks widely utilized in research to evaluate MARL algorithms. Units from the strategy video game StarCraft II engage in confrontations with each other in diver scenarios. The objective is for multiple agents to collaborate in defeating the enemies. There are multiple scenarios, each categorized into difficulty levels such as EASY, HARD, and SuperHARD. We primarily conduct experiments in HARD, and SuperHARD scenarios. Table 3 provides a detailed description of the scenarios we used in our experiments.

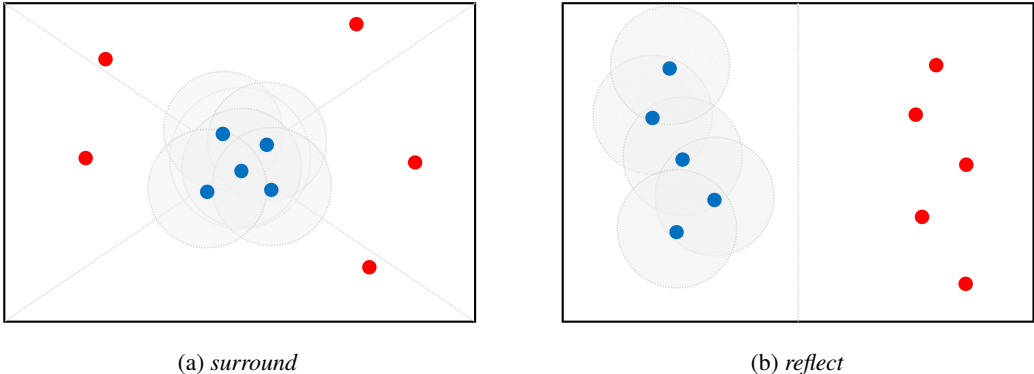
Scenario	Difficulty	Ally Units	Enemy Units	Type
2s_vs_1sc	EASY	2 Stalkers	1 Spine Crawler	micro-trick: alternating fire
3s_vs_3z	EASY	3 Stalkers	3 Zealots	micro-trick: kiting
3s_vs_5z	HARD	3 Stalkers	5 Zealots	micro-trick: kiting
2c_vs_64zg	HARD	2 Colossi	64 Zerglings	micro-trick: positioning
MMM	HARD	1 Medivac 2 Marauders 7 Marines	1 Medivac 2 Marauders 7 Marines	heterogeneous & symmetric
corridor	Super HARD	6 Zealots	24 Zerglings	micro-trick: wall off
6h_vs_8z	Super HARD	6 Hydras	8 Zealots	micro-trick: focus fire
MMM2	Super HARD	1 Medivac 2 Marauders 7 Marines	1 Medivac 3 Marauders 8 Marines	heterogeneous & asymmetric
1o_2r_vs_4r	-	1 Overload 2 Roaches	4 Roaches	communication

Table 3: A detailed description of the SMAC scenario used in the experiment

A.1.2 SMACv2

SMACv2 is proposed to address the shortcomings of SMAC, particularly in terms of its lack of stochasticity and partial observable characteristics (Ellis et al., 2023). Therefore, SMACv2 differs from SMAC in three main aspects.

756 First, the unit composition is randomly determined. In SMAC, the generated units are fixed, whereas
 757 in SMACv2, different types of units are randomly generated based on probabilities. The second
 758 difference lies in the observation probability of agents. In SMAC, when one agent observes an
 759 enemy, other agents can also observe the same enemy simultaneously. In contrast, in SMACv2, if one
 760 agent observes an enemy first, other agents within their observation range may not identify the same
 761 enemy, even if it is present. The last distinction involves adding randomness to the location where
 762 units are spawned. The location where units are spawned is determined by one of two types: one is
 763 *surround* and the other is *reflect*. *surround* entails the creation of units in a formation where allied
 764 units surround enemy units, while *reflect* involves units being spawned in a facing and confronting
 765 manner.



779 Figure 11: Two different types of start positions in SMACv2

781 We conduct experiments in SMACv2 under conditions that encompass all three aforementioned
 782 aspects, and additionally, we fixed the last condition to investigate performance differences in MA²E
 783 based on the spawning locations. To elaborate further, when comparing the results between the
 784 *surround* and *reflect*, we assumed that *surround*, where agent observations already sufficiently
 785 overlap, would render global information less significant. On the other hand, in the case of *reflect*,
 786 agents are spawned at an appropriate distance, allowing for the assumption that through the inference
 787 of global information, they could acquire useful information about agents at long range.

788 Figure 11 provides examples of the *reflect* and *surround* positions. The light gray circles represent
 789 the observation range of allied agents. In the case of *surround*, most areas overlap. In contrast, *reflect*
 790 shows a relatively smaller overlap compared to *surround*. Therefore, we believe that MA²E can better
 791 infer global information in the *reflect* position than in the *surround* position.

793 A.1.3 GOOGLE RESEARCH FOOTBALL (GRF)

794 Google Research Football (GRF) is one of the widely used benchmark in multi-agent reinforcement
 795 learning research. Multiple agents cooperate to play a football game in GRF and GRF offers various
 796 scenarios. We experiment with our model in three different scenarios of GRF.

- 798 • *academy_3_vs_1with_keeper*: The objective is for three allied agents to score goals in a
 799 soccer half-court field against an opposing goalkeeper.
- 800 • *academy_counterattack_easy*: The scenario consists of four allied agents, one opposing
 801 agent, and one goalkeeper. The objective is for the allied agents to execute a counter-attack.
- 802 • *academy_counterattack_hard*: The scenario consists of four allied agents, two opposing
 803 agents, and one goalkeeper. The objective is for the allied agents to execute a counter-attack.

806 A.2 HYPERPARAMETERS

807 The hyperparameters for the baselines in our experiments are as listed in Table 4. The value-based
 808 algorithms include VDN, QMIX, QPLEX, ow-QMIX, cw-QMIX, while the policy-based algorithms
 809 include IPPO, RIIT. Moreover, the hyperparameters for MA²E are listed in Table 5.

Table 4: The hyperparameter settings for the baseline algorithms

Algorithms	Value-based	Policy-based
Optimizer	Adam	Adam
Batch Size	128	64, 32
TD(λ)	0.6	-
Learning Rates	0.001	0.0005, 0.001
Replay Buffer Size	1000	64, 128
ϵ Anneal Steps	100000	100000
Gamma	0.99	

Table 5: The hyperparameter settings for MA²E

Hyperparameters	Value
Batch size	32
Input embedding	24
The number of heads	4
The number of encoder layer	3
The number of decoder layer	2
Steps for fine tuning	500
Pretraining threshold	0.015

B PSEUDOCODE

Algorithm 1 Model with Multi-Agent Masked Auto-Encoder (MA²E) Applied

- 1: Initialized value networks $Q_{\theta_{ind}}$ and $Q_{\theta_{tot}}$ or policy network π_{θ}
 - 2: Initialized MA²E parameters θ_{ma^2e}
 - 3: Prepare replay buffer D
 - 4: **repeat**
 - 5: Run episodes through random policy and store trajectories in Buffer D
 - 6: Update MA²E parameters θ_{ma^2e} using samples in D : $L_{MA^2E} = \frac{1}{nT} \sum_{t=1}^T \sum_{i=1}^n (\tau_t^i - \tilde{\tau}_t^i)^2$
 - 7: **until** When training has been done for a specific number of steps or when the loss is lower than the threshold
 - 8: Reset replay buffer D
 - 9: **repeat**
 - 10: **for** each episode **do**
 - 11: Get initial state s
 - 12: **while** episode is not terminated **do**
 - 13: Sample actions \mathbf{a}_t from Q with ϵ greedy or policy π_{θ}
 - 14: Execute actions and observe reward r_t
 - 15: Store transition $(s_t, \mathbf{o}_t, \mathbf{u}_t, r_t)$ in buffer D
 - 16: **end while**
 - 17: Update value networks $Q_{\theta_{ind}}$ and $Q_{\theta_{tot}}$ or policy network π_{θ}
 - 18: Update MA²E parameters θ_{ma^2e}
 - 19: **end for**
 - 20: **until** reaching maximum total environment steps
-

C ADDITIONAL RESULTS OF MA²E WITH VDN, QPLEX, RIIT AND IPPO ACROSS DIFFERENT ENVIRONMENT

We provide additional experimental results from SMAC, SMAC2, and GRF. All results were tested with 3 random seeds.

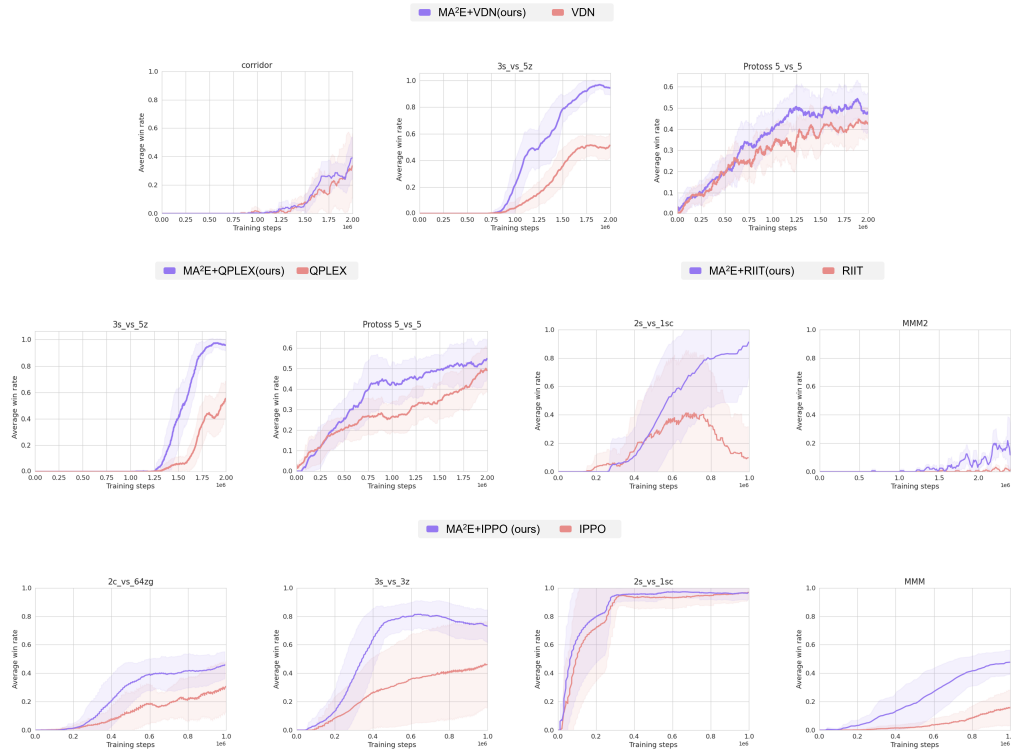


Figure 12: Performance comparison in SMAC and SMACv2 scenarios

D ABLATION STUDIES ON OBSERVATION OVERLAP

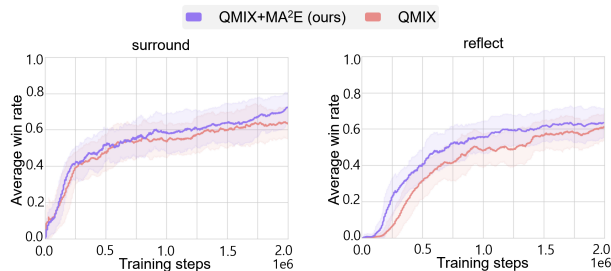


Figure 13: Performance comparison according to the different starting position in SMACv2

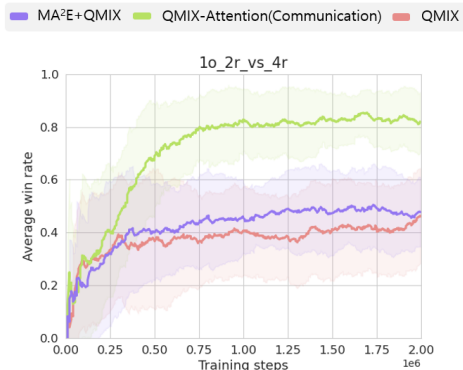
SMACv2 offers different starting positions based on the configurations, with two main strategies for positioning: *surround*, where the units are placed in a position surrounded by the enemy, and *reflect*, where they are positioned head-on in a confrontation with the enemy. The extent of overlap in agent observations varies depending on each position. (For more details, please refer to Appendix A.1.2.) We utilize these two positions to compare the performance of MA²E based on the extent of

918 observation overlap. In both cases, the performance is better than the baseline, but it is particularly
 919 evident from Figure 13 that the performance is significantly better in the *reflect* compared to the
 920 baseline. Therefore, it can be concluded that when the observation range overlaps at a reasonable
 921 proportion, MA²E can better infer global information, as opposed to cases where the overlap is too
 922 extensive.

923
 924
 925 **E LIMITATIONS**
 926
 927



941 (a) A Screenshot of 1o_2r_vs_4r scenario in
 942 SMAC



947 (b) Win rate comparison between using communica-
 948 tion and using MA²E.

949
 950
 951
 952
 953
 954
 955
 956
 957
 958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971

Figure 14: A scenario where MA²E fails to infer global information

947 An agent must use observations as clues to infer the actions or movements of other agents. Therefore,
 948 if there is no information about other agents in the observation of an agent, MA²E cannot effectively
 949 capture the information.

950 For example, in environments where agents are spread far apart due to a large map size, and their
 951 observations do not overlap, MA²E may struggle to infer the information. Figure 14a shows a
 952 screenshot of the 1o_2r_vs_4r scenario, and Figure 14b shows a comparing the win rates when
 953 the model uses communication and when the model uses MA²E. In this map, enemies and allies
 954 spawn at random locations among the four corners (top-left, bottom-left, top-right, bottom-right) of
 955 the map. Allied agents consist of two attacking units and one observation unit, with the observation
 956 unit spawning at the same location as the enemies.

957 If communication is possible, the observation unit can inform the attacking units of the enemy’s
 958 location. On the other hand, since the observation unit is outside the observation range of the attacking
 959 units, they cannot infer information about the observation unit. We compare our proposed model with
 960 the communication method presented by Hu et al. (2021), called QMIX-Attention, to investigate this
 961 situation. As shown in the Figure 14b, the performance of our proposed model is relatively lower
 962 compared to the case where communication is used. Nevertheless, despite that, it can be confirmed
 963 that using MA²E does not make lower performance compared to the default model.

964 We also compare with communication methods Commnet (Sukhbaatar et al., 2016), G2ANet (Liu
 965 et al., 2020), MAIC (Yuan et al., 2022), and QMIX-att (Hu et al., 2021). The results are shown in
 966 Figure 15. In 3s_vs_5z, appropriately overlapping observations allow MA²E to effectively infer
 967 information, and the use of attention layers reduces the space size. Conversely, communication based
 968 methods increase the space size and require communication overhead, leading to slower convergence
 969 and sample inefficiency. However, in 1o_2r_vs_4r, MA²E struggles to perform effectively due
 970 to a very little observation overlap among agents, while communication methods (QMIX-att and
 971 Commnet) can still fully share the information hence achieves better performance. Thus, we note that
 the need for communication vary depending on the environment.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

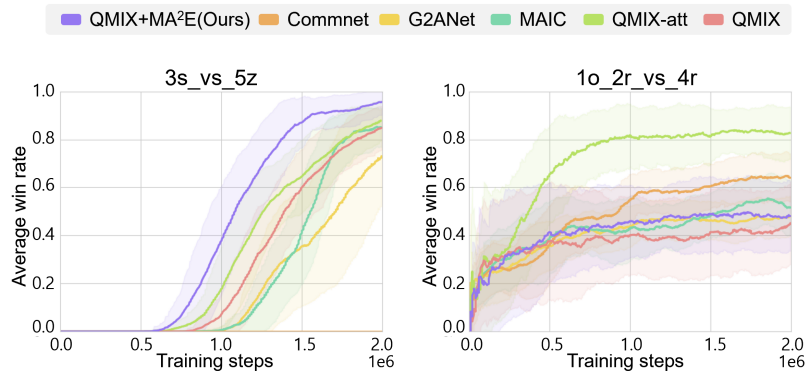
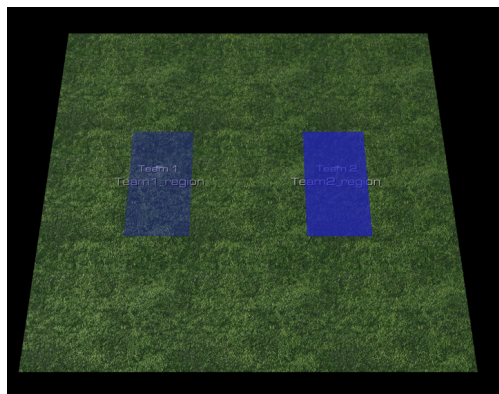


Figure 15: Performance comparison with communication based MARL algorithms in two SMAC scenarios

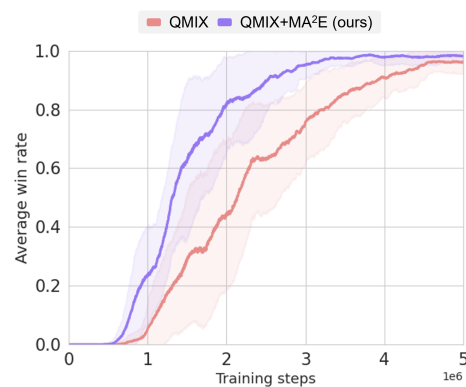
F IS MA²E PERFORMING WELL BECAUSE IT IS OVERFITTED TO THE ENVIRONMENT?

MA²E learns the dynamics of the environment during the training process and captures global information through its ability to infer masked areas. Therefore, from one perspective, it can be considered that MA²E performs well due to overfitting to the environment. Specifically, in the SMAC environment, there is a significant performance difference between algorithms applying MA²E and the baseline, but in the SMACv2 environment, this difference is not as pronounced. This suggests that MA²E may be overfitting to the environment because SMAC is more deterministic compared to SMACv2, making it easier for the model to overfit.

However, the issue of not performing well in SMACv2 is a common difficulty not only for MA²E but also for existing MARL algorithms. In reflecting the stochastic characteristic of SMACv2, the random properties can generate scenarios that are unwinnable from the start or turn each episode into a completely different task, requiring a multi-task approach. Therefore, it is difficult to argue that MA²E is overfitting just because it does not perform well in SMACv2.



(a) The spawn locations of enemy and ally agents in the modified 3s_vs_5z scenario.



(b) Win rate comparison between using MA²E and without MA²E case

Figure 16: The experimental results in the environment with added stochastic characteristics

Nonetheless, to prove that MA²E is not overfitting to the environment, we created a non-deterministic scenario by adding random characteristics to the 3s_vs_5z scenario of SMAC and used it for experimentation. In the original 3s_vs_5z scenario, the positions where enemy and ally agents

1026 spawn are fixed. In contrast, in our modified scenario, the enemy and ally agents spawn at random
1027 locations within the blue area shown on the Figure 16a. Thus, a stochastic property is added where
1028 the spawn positions of the enemy and ally agents change in each episode.

1029 The experimental results are shown on the Figure 16b. As a result, it can be seen that the model
1030 using MA²E significantly outperforms the model without it. Therefore, it can be concluded that
1031 MA²E does not overfit to the environment but appropriately infers global information according to
1032 the situation.

1033

1034

1035 G BROADER IMPACT

1036

1037 Our study introduces Multi-Agent Masked Auto-Encoder (MA²E) to enhance decision-making
1038 in multi-agent systems. By addressing partial observability through a masking perspective, our
1039 approach has broad applications in real-world scenarios such as military operations, autonomous
1040 driving, traffic systems, and robotics, promising improved decision-making across diverse multi-agent
1041 environments. While explicit communication can selectively send the messages, MA²E is trained
1042 with full information. Hence the proposed method may cause privacy or security issues in some real
1043 world applications, which would be mitigated by anonymizing data.

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

H ADDITIONAL EXPERIMENT AND RESULTS FOR REBUTTAL

H.1 RESULTS IN SMAC

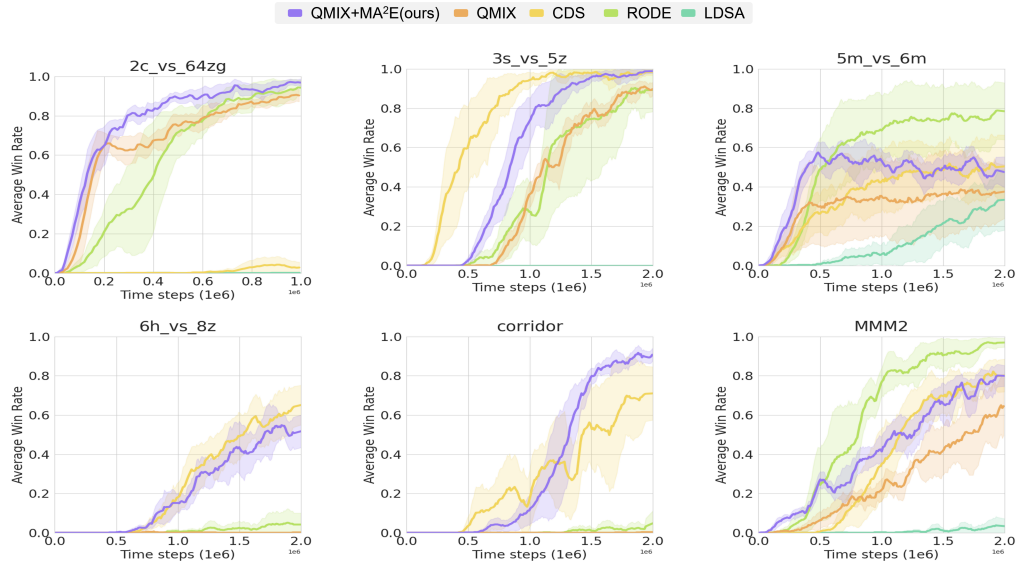


Figure 17: Additional experimental results comparing QMIX+MA²E with LDSA, CDS, and RODE in SMAC. While each baselines suffers in several scenarios, QMIX+MA²E consistently performs well across different scenarios.

Table 6: Comparison of QMIX+MA²E with the baselines, corresponding to Figure 17. The values in the table represent average win rates with standard deviations. QMIX+MA²E outperforms baselines on average across different scenarios.

Scenario	Difficulty	Steps	QMIX+MA ² E (ours)	QMIX	LDSA	RODE	CDS
3s_vs_5z	HARD	2M	0.99±0.01	0.89±0.03	0.0±0.0	0.89±0.11	0.98±0.01
2c_vs_64zg	HARD	1M	0.97±0.02	0.9±0.03	0.01±0.0	0.94±0.04	0.03±0.03
5m_vs_6m	HARD	2M	0.48±0.08	0.37±0.14	0.33±0.0	0.78±0.14	0.5±0.16
corridor	SuperHARD	2M	0.9±0.03	0.0±0.0	0.0±0.0	0.04±0.06	0.71±0.13
6h_vs_8z	SuperHARD	2M	0.52±0.08	0.0±0.0	0.0±0.0	0.04±0.05	0.64±0.1
MMM2	SuperHARD	2M	0.81±0.05	0.64±0.15	0.03±0.04	0.96±0.02	0.8±0.08
Avg.	-	-	0.78	0.47	0.06	0.61	0.61

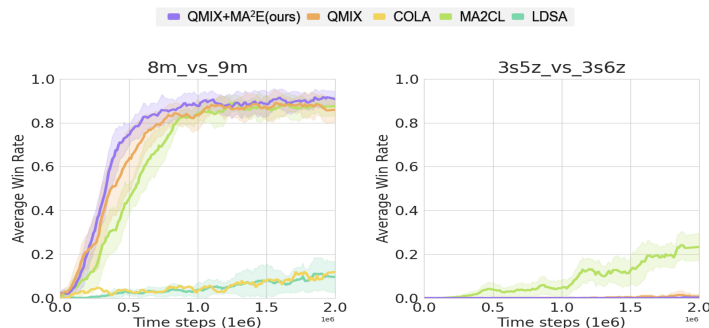


Figure 18: Experimental results in additional SMAC scenarios.

H.2 RESULTS IN SMACv2

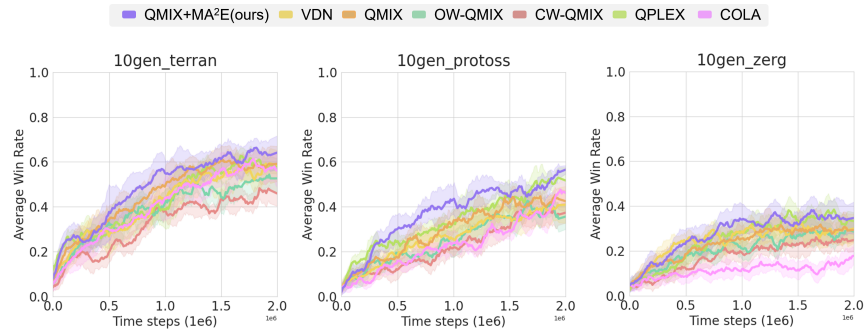


Figure 19: Additional experimental results in SMACv2, with 10gen_zerg is added to the experiment from the original manuscript. QMIX+MA²E outperforms the baselines across different scenarios.

H.3 LEARNING CURVE OF TABLE 2

We provide learning curves of each algorithm in Table 2. For ease of reference, we include here Table 7, which is the same as Table 2. As shown in Figure 20, QMIX+MA²E outperforms the baselines in most scenarios.

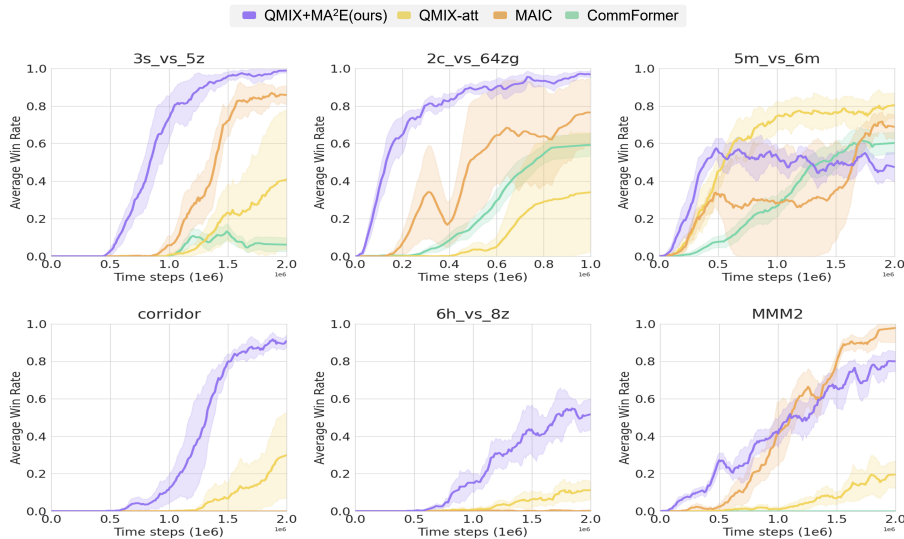


Figure 20: Learning curves of each algorithm in Table 2.

Table 7: Comparison of MA²E with communication-based methods

Scenario	Steps	QMIX+MA ² E (ours)	MAIC	QMIX-att	CommFormer
3s_vs_5z	2M	0.99 ± 0.01	0.86 ± 0.05	0.41 ± 0.37	0.06 ± 0.04
2c_vs_64zg	1M	0.97 ± 0.03	0.76 ± 0.18	0.34 ± 0.31	0.59 ± 0.06
5m_vs_6m	2M	0.48 ± 0.08	0.69 ± 0.07	0.81 ± 0.04	0.61 ± 0.05
corridor	2M	0.91 ± 0.10	0.0 ± 0.0	0.30 ± 0.23	0.0 ± 0.0
6h_vs_8z	2M	0.53 ± 0.08	0.0 ± 0.0	0.14 ± 0.04	0.0 ± 0.0
MMM2	2M	0.79 ± 0.08	0.97 ± 0.07	0.20 ± 0.08	0.0 ± 0.0

H.4 ABLATION STUDY

We provide additional ablation study in `2c_vs_64zg` scenario. For ease of comparison, we include here Figure 22, which is the same as Figure 10. We can draw a similar conclusion that: using agent-level masking, an appropriate number of trajectories, positional encoding considering both agent and time, and both pretraining and fine-tuning yield better results.

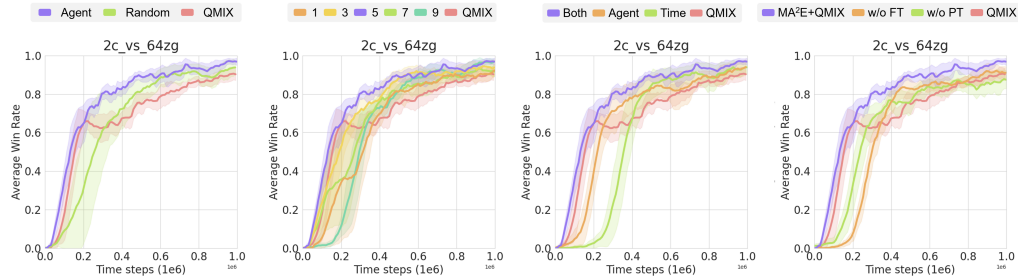


Figure 21: Additional ablation study in `2c_vs_64zg` scenario. Along with Figure 22, we can draw a similar conclusion.

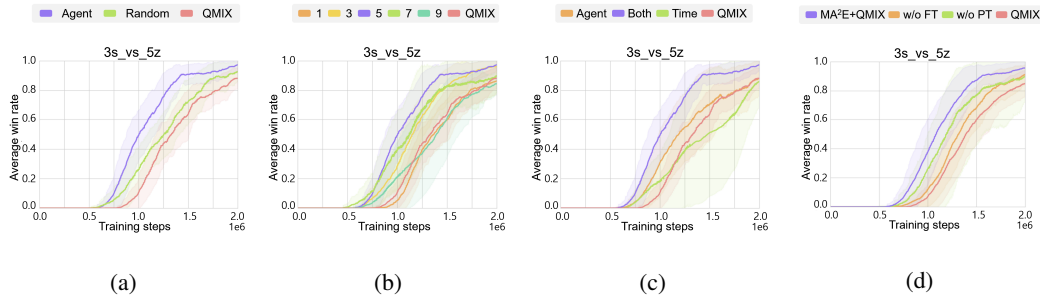


Figure 22: Performance comparison based on hyper parameter settings or strategies in `3s_vs_5z` scenario (same as Figure 10). **(a)**: Performance comparison based on masking strategies. The agent-based masking strategy applied model outperforms both the random masking strategy applied model and the default model. **(b)**: Performance comparison of MA^2E based on the number of trajectories. **(c)**: Performance comparison of MA^2E based on the positional encoding strategies. When considering both agent and time, it shows the best performance. **(d)**: Performance comparison of MA^2E with or without pre-training and fine-tuning.

H.5 COMPARISON OF TRAINING COST AND MODEL SIZE WITH AND WITHOUT USING MA²E

3s_vs_5z	MA2E+QMIX	QMIX	Difference
Time for pretraining	0.77h	-	+0.77h
Time to convergence	17.05h (faster)	18.22h	-1.17h (93%)
Steps to convergence	2.03M (faster)	3.52M	-1.49M (57%)
Steps per second	33.07	53.66	-20.59
The number of parameters	0.43M	0.12M	0.31M

Table 8: Comparison of training and execution time and the number of parameters for each method in the SMAC `3s_vs_5z` scenario. The values in parenthesis represent the relative ratio between MA2E+QMIX and QMIX.

2c_vs_64zg	MA2E+QMIX	QMIX	Difference
Time for pretraining	0.84h	-	+0.84h
Time to convergence	7.97h (faster)	8.76h	-0.79h (90%)
Steps to convergence	0.639M (faster)	1.021M	-0.382M (63%)
Steps per second	22	32	-10
The number of parameters	0.98M	0.65M	0.33M

Table 9: Comparison of training and execution time and the number of parameters for each method in the SMAC `2c_vs_64zg` scenario. The values in parenthesis represent the relative ratio between MA2E+QMIX and QMIX.

H.6 RESULTS IN SMAC UNDER THE SAME HYPERPARAMETER SETTING AND PYMARL2

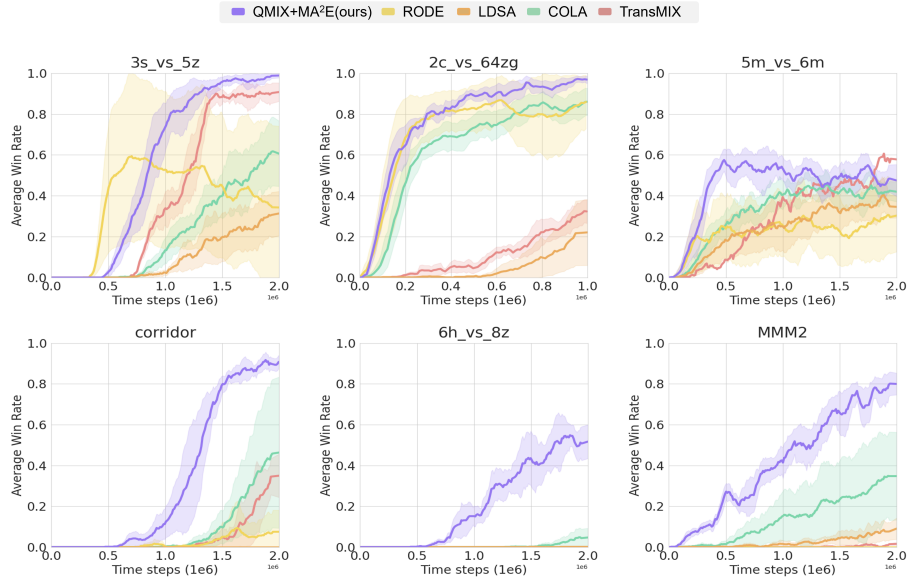


Figure 23: Additional experimental results in SMAC. All the baselines are evaluated under the same hyperparameter settings as MA²E and pymarl2. Again, QMIX+MA²E consistently performs well across different scenarios.

Table 10: Comparison of QMIX+MA²E with the baselines, corresponding to Figure 23. The values in the table represent average win rates with standard deviations. QMIX+MA²E outperforms baselines on average across different scenarios.

Scenario	Difficulty	Steps	QMIX+MA ² E (ours)	COLA	RODE	LDSA	TransMIX
3s_vs_5z	HARD	2M	0.99±0.01	0.60±0.16	0.34±0.4	0.31±0.10	0.91±0.04
2c_vs_64zg	HARD	1M	0.97±0.02	0.86±0.07	0.86±0.13	0.22±0.16	0.32±0.05
5m_vs_6m	HARD	2M	0.48±0.08	0.37±0.14	0.42±0.05	0.30±0.17	0.58±0.01
corridor	SuperHARD	2M	0.9±0.03	0.46±0.37	0.08±0.11	0.0±0.0	0.35±0.11
6h_vs_8z	SuperHARD	2M	0.52±0.08	0.05±0.05	0.0±0.0	0.0±0.0	0.0±0.0
MMM2	SuperHARD	2M	0.81±0.05	0.35±0.22	0.0±0.0	0.09±0.05	0.02±0.00
Avg.	-	-	0.78	0.46	0.26	0.16	0.36

H.7 COMPARISON WITH THE METHODS EMPLOYING TRANSFORMER IN THE AGENT NETWORK

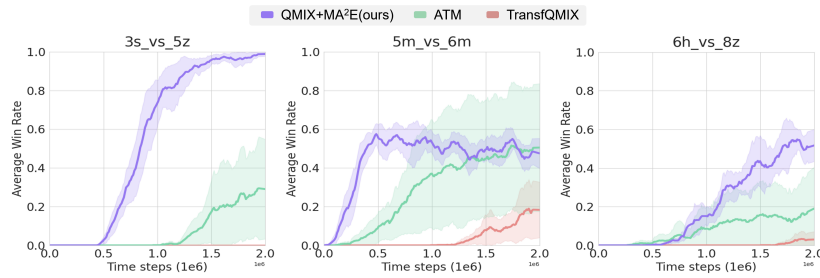


Figure 24: Additional experimental results comparing MA²E with baselines that employ the Transformer architecture in agent networks.

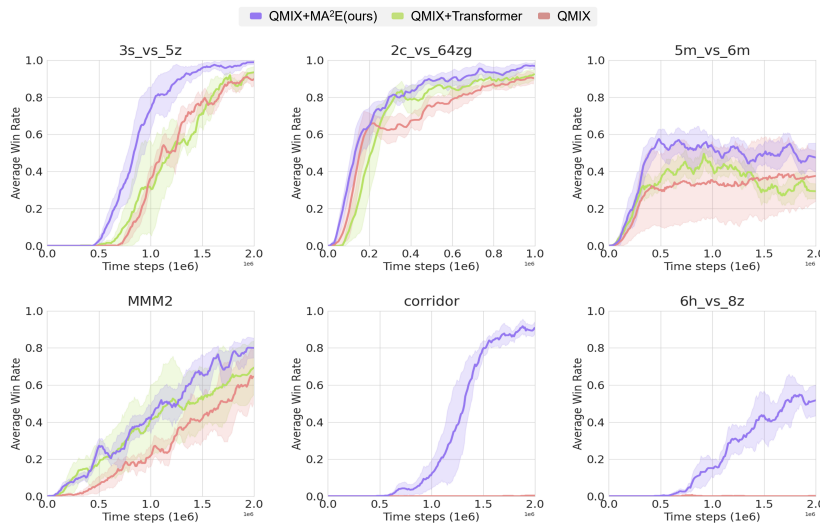


Figure 25: Additional experimental results on MA²E without masking and its associated loss while retaining the Transformer architecture in the agent network (labeled as QMIX+Transformer).

H.8 FURTHER EXPERIMENTAL RESULTS FOR RIIT

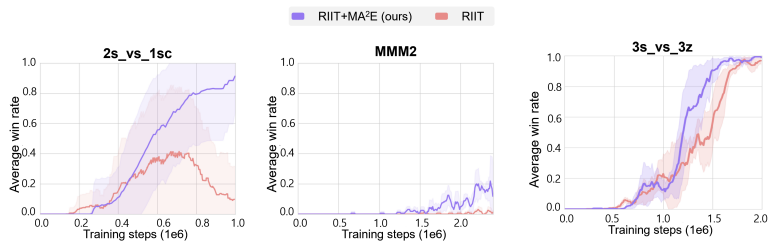


Figure 26: Additional experimental results employing MA²E on RIIT.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457