

---

# SMACv2: An Improved Benchmark for Cooperative Multi-Agent Reinforcement Learning

---

Benjamin Ellis<sup>1\*</sup> Jonathan Cook<sup>1</sup> Skander Moalla<sup>1</sup> Mikayel Samvelyan<sup>2 3</sup>  
Mingfei Sun<sup>1</sup> Anuj Mahajan<sup>1</sup> Jakob N. Foerster<sup>1</sup> Shimon Whiteson<sup>1</sup>

<sup>1</sup>University of Oxford <sup>2</sup>University College London <sup>3</sup>Meta AI

## 1 Links

- SMACv2: <https://github.com/oxwhirl/smacv2>. MIT License.
- MAPPO implementation: <https://github.com/benellis3/mappo>.
- QMIX implementation: <https://github.com/benellis3/pymarl2>. Apache-2 License.

## 2 Further Background

### 2.1 StarCraft

StarCraft II is a real-time strategy game featuring 3 different races, Protoss, Terran and Zerg, with different properties and associated strategies. The objective is to build an army powerful enough to destroy the enemy’s base. When battling two armies, players must ensure army units are acting optimally. This is called *micromanagement*. An important micromanagement strategy is *focus firing*, which is ordering all allied units to jointly target the enemies one by one to ensure that damage taken is minimised.

Another important strategy is *kiting*, where units flee from the enemy and then pick them off one by one as they chase.

### 2.2 QMIX

QMIX can be thought of as an extension of DQN[7] to the Dec-POMDP setting. The joint optimal action is found by forcing the joint  $Q$  to adhere to the individual global max (IGM) principle[11], which states that the joint action can be found by maximising individual agents’  $Q_i$  functions:

$$\arg \max_{\mathbf{a}} Q(s, \boldsymbol{\tau}, \mathbf{a}) = \begin{cases} \arg \max_{a_1} Q_1(\tau_1, a_1) \\ \arg \max_{a_2} Q_2(\tau_2, a_2) \\ \dots \\ \arg \max_{a_n} Q_n(\tau_n, a_n) \end{cases}$$

This central  $Q$  is trained to regress to a target  $r + \gamma \hat{Q}(s, \boldsymbol{\tau}, \mathbf{a})$  where  $\hat{Q}$  is a target network that is updated slowly. The central  $Q$  estimate is computed by a mixing network, whose weights are conditioned on the state, which takes as input the utility function  $Q_i$  of the agents. The weights of the mixing network are restricted to be positive, which enforces the IGM principle[11] by ensuring the central  $Q$  is monotonic in each  $Q_i$ .

---

\*Correspondence to [benellis@robots.ox.ac.uk](mailto:benellis@robots.ox.ac.uk).

## 24 2.3 Independent and Multi-agent PPO

25 Proximal Policy Optimisation (PPO) is a method initially developed for single-agent reinforcement  
26 learning which aims to address performance collapse in policy gradient methods. It does this by  
27 heuristically bounding the ratio of action probabilities between the old and new policies. To this end,  
28 it optimises the below objective function.

$$\mathbb{E}_{s \sim d^\pi, a \sim \pi} \left[ \min \left( \frac{\tilde{\pi}(a|s)}{\pi(a|s)} A^\pi(s, a), \text{clip} \left( \frac{\tilde{\pi}(a|s)}{\pi(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A^\pi(s, a) \right) \right]$$

29 where  $\text{clip}(t, a, b)$  is a function that outputs  $a$  if  $t < a$ ,  $b$  if  $t > b$  and  $t$  otherwise.

30 Extending this to the multi-agent setting can easily be done in two ways. The first is to use independent  
31 learning, where each agent treats the others as part of the environment and learns a critic using its  
32 observation and action history. The second is to use a centralised critic conditioned on the central  
33 state. This is called multi-agent PPO. Both independent PPO (IPPO) and multi-agent PPO (MAPPO)  
34 have demonstrated strong performance on SMAC [16, 2]. Note that we do not apply the observation  
35 and state changes suggested by Yu et al. [16] anywhere in this paper. This is because these changes  
36 were not implemented as a wrapper on top of SMAC, but instead by modifying SMAC directly,  
37 leading to the environment code becoming unmanageably complicated.

## 38 3 Experimental Details

39 In this section we describe extra details of the experiments run as part of the paper.

### 40 3.1 Stochasticity

41 This section describes details of the experiments run in Section 5.1 in the main paper. Both the  
42 closed-loop and open-loop algorithms were based on the MAPPO implementation from Sun et al.  
43 [13]. The code used for these experiments can be found here. Both the open-loop and closed-loop  
44 algorithms use the same neural network architecture as in [13]. Both were also provided with access  
45 to the available actions mask of the environment and conditioned the critic on the state. We used the  
46 hyperparameters from Sun et al. [13], with the exception of the actor’s learning rate, which we set  
47 to 0.0005 and decayed linearly throughout training. This is because learning rate decay has been  
48 found to be important for bounding PPO’s policy updates [12]. Full hyperparameters can be found in  
49 Table 5. Full results are shown in Figure 1.

### 50 3.2 SMAC Feature Inferrability & Relevance

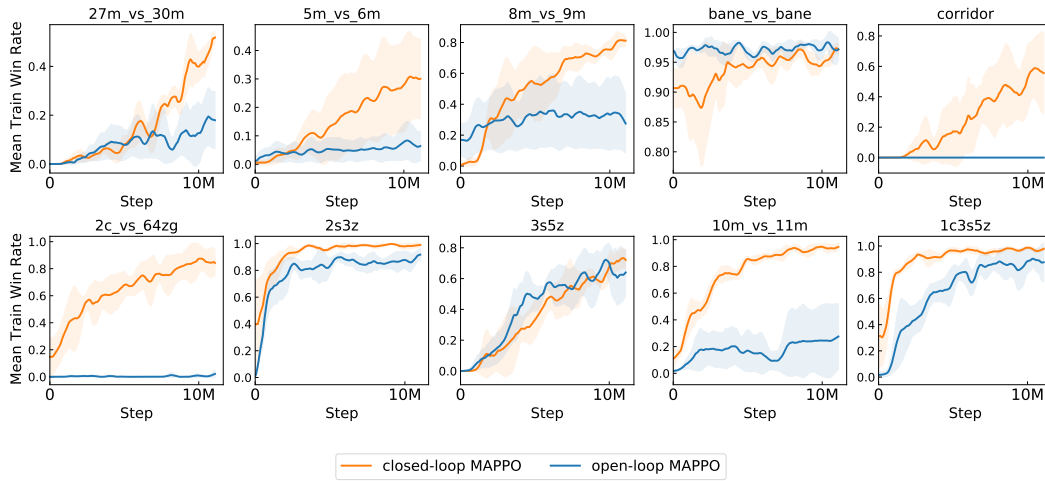
51 This section describes details of the experiments run in Section 5.2 and Section 7.3 in the main paper.  
52 The code used for these experiments can be found here.

53 For each scenario, we had 3 QMIX policies trained using the implementation and hyperparameters  
54 from Hu et al. [3], each with a different network and SMAC initialization. Training results for the  
55 policies for SMAC are shown in Figure 2. Each policy constitutes a seed and is used to collect a  
56 dataset of episodes for the feature-relevance experiment on the scenario. We call these policies *expert*  
57 *policies*. QMIX hyperparameters used are given in Table 4. A dataset consists of two folds: 8192  
58 episodes used for training and 4096 episodes used for evaluation.

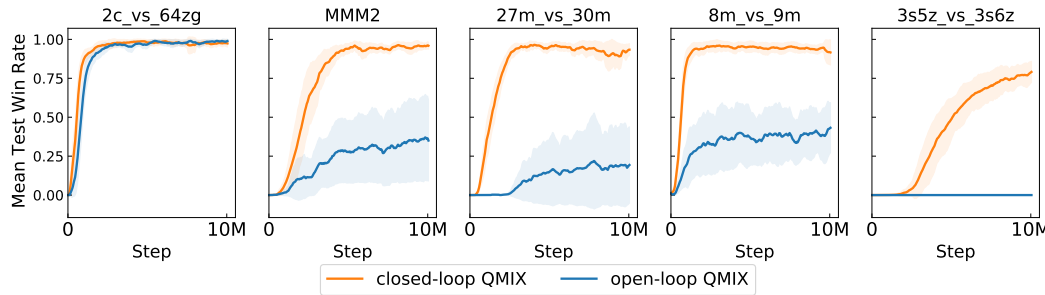
59 For a given mask, the experiment then consists of training a new QMIX network, termed *regression*  
60 *network*, whose input is trajectories with observations and states masked according to the mask, and  
61 whose task is to predict the  $Q$ -values output by the expert policy on the unmasked trajectories.

62 The *nothing* mask does not apply any effect on the observations and states in a trajectory. Otherwise,  
63 a mask, say ‘health (ally)’, masks the health feature of every ally in the observation of each agent  
64 (except its own health) and masks the health feature of all the units controlled by QMIX in the state.  
65 Table 1 shows the feature sets zeroed out for different masks.

66 The regression network has the same architecture as the expert network. We use the mean squared  
67 error (MSE) as a loss function and optimise the network via Adam with a batch size of 512 episodes  
68 and learning rate equal to 0.005. (The other Adam hyper-parameters are the default PyTorch values.)



(a) MAPPO open-loop and closed-loop results



(b) QMIX open-loop and closed-loop results

Figure 1: Plot of selected SMAC scenarios treated as an open-loop planning problem by limiting the observation to the current timestep and agent ID. Plots show the mean win rate and standard deviation across 3 training seeds for MAPPO and QMIX.

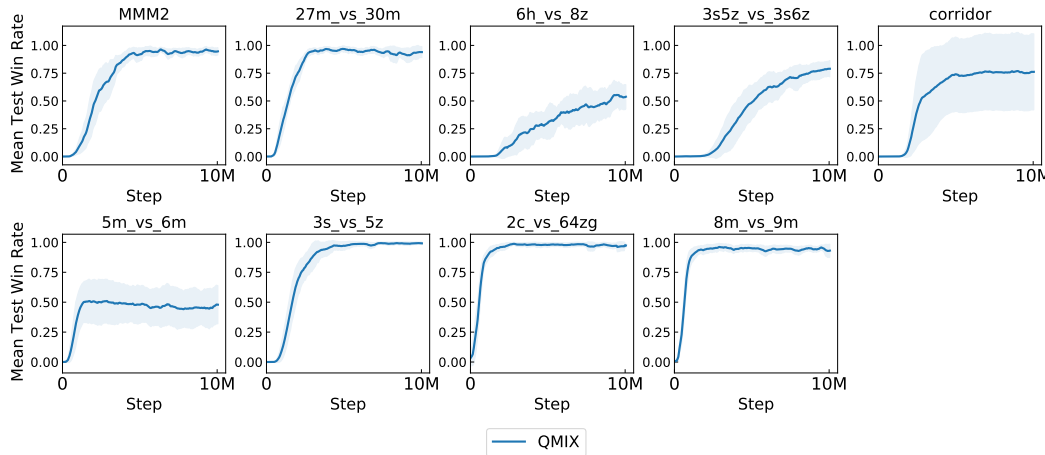


Figure 2: QMIX training results on SMAC

Table 1: Masked features for each setting of the feature quality experiment. ✓ means a feature is masked and ✗ means a feature is not masked.

Mask	Ally						Enemy				
	Health	Shield	$x$	$y$	Distance	Actions	Health	Shield	$x$	$y$	Distance
<i>everything</i>	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
<i>nothing</i>	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
<i>health (ally)</i>	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
<i>shield (ally)</i>	✗	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗
<i>distance (ally)</i>	✗	✗	✓	✓	✓	✗	✗	✗	✗	✗	✗
<i>health and shield (ally)</i>	✓	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗
<i>actions only</i>	✗	✗	✗	✗	✗	✓	✗	✗	✗	✗	✗
<i>all except actions</i>	✓	✓	✓	✓	✓	✗	✗	✗	✗	✗	✗
<i>ally all</i>	✓	✓	✓	✓	✓	✓	✗	✗	✗	✗	✗
<i>health (enemy)</i>	✗	✗	✗	✗	✗	✗	✓	✗	✗	✗	✗
<i>shield (enemy)</i>	✗	✗	✗	✗	✗	✗	✗	✓	✗	✗	✗
<i>distance (enemy)</i>	✗	✗	✗	✗	✗	✗	✗	✗	✓	✓	✓
<i>enemy all</i>	✗	✗	✗	✗	✗	✗	✓	✓	✓	✓	✓

69 Training is performed with early stopping according to the validation fold with an evaluation every 5  
70 epochs and a patience of 10 tries (i.e. training is stopped when the MSE on the validation fold does  
71 not decrease in 50 consecutive epochs and performance at the best epoch is retained). Models for all  
72 masks and scenarios hit early stopping and trained for about 200 epochs on average, except a few  
73 ones which trained to a limit of 500 epochs.

74 All hyper-parameters were tuned on a few different scenarios from both SMAC and SMACv2 to  
75 minimise validation MSE using datasets and expert policies not used for the reported results. The  
76 sizes of the training and validation folds are respectively 1.6 times and 0.8 times the size of the QMIX  
77 replay buffer (5000 episodes) and have been chosen to be large enough to minimize validation MSE  
78 while keeping experiments practical.

79 Figure 3 shows the full results of the feature quality experiments for SMAC and Figure 4 shows the  
80 full results for SMACv2. Table 3 shows summary results for the *everything* and *nothing* masks for  
81 SMACv2 and Table 2 shows the same data for SMAC.

82 Experiments for this section were conducted on 80-core CPU machines with NVIDIA GeForce RTX  
83 2080 Ti or Tesla V100-SXM2-16GB GPUs. A single (scenario, seed, mask) combination took around  
84 1 hour to train. In total, the experiments in this section took about 1500 GPU hours.

### 85 3.3 SMACv2 Runs

86 Here we describe the hyperparameter and training procedure used in section 7.1. As in previous  
87 experiments, we used the implementation by Hu et al. [3] for QMIX and by Sun et al. [13] for  
88 MAPPO. These implementations have both achieved very strong results on SMAC. We therefore  
89 tuned the hyperparameters by taking deviations from these for a few key parameters. For MAPPO,  
90 we tuned the actor’s learning rate and the clipping range. Additionally, we added linear learning rate  
91 decay to the MAPPO implementation because this has been shown to be important to bound the policy  
92 update[12]. For QMIX, we tuned the  $\epsilon$ -annealing time and  $\lambda$  in the eligibility trace  $Q(\lambda)$ . All other  
93 hyperparameters, including neural network architecture, were unchanged from the implementations  
94 mentioned previously.

95 The MAPPO code for these experiments can be found here, and the QMIX code here. The QMIX  
96 code is distributed under the Apache license, and the MAPPO code under the MIT license. The only  
97 differences in these branches and the implementations used for previous experiments on SMAC are  
98 important to changing the environment from SMAC to SMACv2.

99 The open-loop policy is identical to MAPPO, except the policy receives as input only the timestamp  
100 and agent ID, as in Section 5.1.

101 We tuned hyperparameters by running each set of hyperparameters on all scenarios and then choosing  
102 the best results. The grids of hyperparameters for QMIX and MAPPO can be found in Table 7  
103 and 8. **For IPPO, we used mostly the same hyperparameters as MAPPO. For QPLEX we used the**  
104 **hyperparameters from [15].** All algorithms were trained for 10M environment steps with evaluations

Table 2: Mean Q values in the different feature quality experiments on SMAC scenarios

Map	Mask	$\bar{Q}$	$\epsilon_{\text{rme}}$	$\frac{\sigma_{\text{me}}}{\bar{Q}}$	$\epsilon_{\text{abs}}$	$\frac{\text{rank} - \text{rank}_{\text{min}}}{Q}$
corridor	everything	7.0 ± 0.08	0.55 ± 0.08	0.078 ± 0.002	0.39 ± 0.05	0.042 ± 0.001
	ally_all_except_actions		0.34 ± 0.09	0.047 ± 0.005	0.21 ± 0.04	0.01 ± 0.008
	ally_health_and_shield		0.29 ± 0.07	0.041 ± 0.003	0.19 ± 0.03	0.05 ± 0.006
	ally_health		0.29 ± 0.05	0.0413 ± 0.0007	0.19 ± 0.02	0.005 ± 0.004
	ally_distance		0.27 ± 0.03	0.038 ± 0.002	0.174 ± 0.009	0.001 ± 0.001
	ally_shield		0.26 ± 0.02	0.037 ± 0.003	0.166 ± 0.005	0.0004 ± 0.0003
	ally_all		0.4 ± 0.2	0.06 ± 0.01	0.28 ± 0.09	0.02 ± 0.02
	enemy_health_and_shield		0.29 ± 0.05	0.0414 ± 0.0005	0.2 ± 0.03	0.005 ± 0.004
	enemy_health		0.3 ± 0.04	0.043 ± 0.003	0.21 ± 0.02	0.006 ± 0.001
	enemy_distance		0.32 ± 0.04	0.046 ± 0.002	0.203 ± 0.009	0.009 ± 0.001
	enemy_shield		0.25 ± 0.03	0.036 ± 0.003	0.161 ± 0.005	-0.0001 ± 0.0008
	enemy_all		0.42 ± 0.09	0.059 ± 0.003	0.29 ± 0.06	0.022 ± 0.006
	nothing		0.25 ± 0.02	0.037 ± 0.003	0.164 ± 0.006	
	3a5z_vs_3a6z	everything	6.97 ± 0.01	0.61 ± 0.01	0.087 ± 0.001	0.418 ± 0.006
ally_all_except_actions			0.34 ± 0.01	0.049 ± 0.002	0.193 ± 0.01	0.0083 ± 0.0009
ally_health_and_shield			0.326 ± 0.01	0.047 ± 0.002	0.182 ± 0.009	0.006 ± 0.001
ally_health			0.333 ± 0.006	0.0479 ± 0.0008	0.188 ± 0.001	0.007 ± 0.001
ally_distance			0.29 ± 0.02	0.041 ± 0.003	0.16 ± 0.01	-0.0 ± 0.003
ally_shield			0.299 ± 0.007	0.043 ± 0.001	0.167 ± 0.004	0.0018 ± 0.0004
ally_all			0.53 ± 0.02	0.076 ± 0.003	0.35 ± 0.02	0.035 ± 0.004
enemy_health_and_shield			0.38 ± 0.02	0.055 ± 0.003	0.25 ± 0.02	0.014 ± 0.002
enemy_health			0.36 ± 0.01	0.051 ± 0.002	0.22 ± 0.01	0.01 ± 0.001
enemy_distance			0.32 ± 0.01	0.046 ± 0.002	0.18 ± 0.008	0.005 ± 0.001
enemy_shield			0.302 ± 0.004	0.0434 ± 0.0007	0.188 ± 0.006	0.0023 ± 0.0008
enemy_all			0.5 ± 0.04	0.071 ± 0.007	0.33 ± 0.04	0.03 ± 0.006
nothing			0.286 ± 0.008	0.041 ± 0.001	0.161 ± 0.007	
5m_vs_6m		everything	8.1 ± 0.08	1.71 ± 0.08	0.21 ± 0.02	1.24 ± 0.09
	ally_all_except_actions		1.17 ± 0.04	0.15 ± 0.02	0.76 ± 0.06	0.007 ± 0.004
	ally_health_and_shield		1.13 ± 0.05	0.14 ± 0.02	0.71 ± 0.05	0.002 ± 0.002
	ally_health		1.13 ± 0.05	0.14 ± 0.02	0.72 ± 0.05	0.001 ± 0.006
	ally_distance		1.15 ± 0.05	0.15 ± 0.02	0.73 ± 0.05	0.004 ± 0.004
	ally_shield		1.12 ± 0.07	0.14 ± 0.02	0.71 ± 0.05	0.001 ± 0.002
	ally_all		1.24 ± 0.06	0.16 ± 0.02	0.83 ± 0.09	0.014 ± 0.007
	enemy_health_and_shield		1.28 ± 0.03	0.16 ± 0.02	0.81 ± 0.03	0.021 ± 0.001
	enemy_health		1.27 ± 0.05	0.16 ± 0.02	0.82 ± 0.06	0.019 ± 0.004
	enemy_distance		1.22 ± 0.02	0.15 ± 0.02	0.82 ± 0.02	0.015 ± 0.007
	enemy_shield		1.11 ± 0.06	0.14 ± 0.02	0.7 ± 0.06	-0.0 ± 0.003
	enemy_all		1.44 ± 0.06	0.18 ± 0.02	0.99 ± 0.07	0.04 ± 0.009
	nothing		1.11 ± 0.05	0.14 ± 0.02	0.71 ± 0.04	
	8m_vs_9m	everything	12.3 ± 0.07	0.75 ± 0.07	0.061 ± 0.006	0.46 ± 0.06
ally_all_except_actions			0.6 ± 0.04	0.049 ± 0.004	0.32 ± 0.04	0.013 ± 0.003
ally_health_and_shield			0.53 ± 0.05	0.043 ± 0.005	0.27 ± 0.04	0.007 ± 0.002
ally_health			0.56 ± 0.03	0.045 ± 0.002	0.29 ± 0.02	0.01 ± 0.004
ally_distance			0.49 ± 0.06	0.04 ± 0.005	0.27 ± 0.05	0.0039 ± 0.0006
ally_shield			0.47 ± 0.02	0.038 ± 0.002	0.27 ± 0.02	0.002 ± 0.003
ally_all			0.73 ± 0.07	0.06 ± 0.007	0.43 ± 0.06	0.023 ± 0.004
enemy_health_and_shield			0.55 ± 0.04	0.045 ± 0.004	0.31 ± 0.03	0.009 ± 0.002
enemy_health			0.55 ± 0.05	0.045 ± 0.004	0.32 ± 0.03	0.009 ± 0.001
enemy_distance			0.46 ± 0.07	0.038 ± 0.006	0.27 ± 0.04	0.002 ± 0.002
enemy_shield			0.42 ± 0.05	0.034 ± 0.005	0.23 ± 0.04	-0.001 ± 0.002
enemy_all			0.64 ± 0.06	0.052 ± 0.006	0.38 ± 0.05	0.016 ± 0.003
nothing			0.44 ± 0.07	0.036 ± 0.006	0.24 ± 0.04	
27m_vs_30m		everything	9.6 ± 0.04	0.49 ± 0.04	0.051 ± 0.006	0.33 ± 0.03
	ally_all_except_actions		0.47 ± 0.05	0.049 ± 0.007	0.28 ± 0.04	0.014 ± 0.003
	ally_health_and_shield		0.44 ± 0.07	0.046 ± 0.009	0.26 ± 0.05	0.01 ± 0.003
	ally_health		0.45 ± 0.06	0.047 ± 0.007	0.26 ± 0.04	0.011 ± 0.002
	ally_distance		0.35 ± 0.06	0.037 ± 0.008	0.21 ± 0.04	0.001 ± 0.001
	ally_shield		0.35 ± 0.02	0.036 ± 0.004	0.2 ± 0.02	0.001 ± 0.003
	ally_all		0.5 ± 0.06	0.052 ± 0.008	0.32 ± 0.05	0.016 ± 0.003
	enemy_health_and_shield		0.46 ± 0.04	0.048 ± 0.006	0.31 ± 0.03	0.012 ± 0.004
	enemy_health		0.5 ± 0.05	0.052 ± 0.007	0.34 ± 0.04	0.0166 ± 0.0063
	enemy_distance		0.41 ± 0.03	0.042 ± 0.005	0.25 ± 0.03	0.007 ± 0.003
	enemy_shield		0.35 ± 0.06	0.036 ± 0.007	0.21 ± 0.03	0.0009 ± 0.0008
	enemy_all		0.53 ± 0.02	0.055 ± 0.004	0.37 ± 0.02	0.02 ± 0.002
	nothing		0.34 ± 0.05	0.036 ± 0.006	0.21 ± 0.04	
	2c_vs_64zg	everything	7.72 ± 0.02	0.56 ± 0.02	0.072 ± 0.002	0.42 ± 0.01
ally_all_except_actions			0.15 ± 0.04	0.02 ± 0.005	0.08 ± 0.02	0.0018 ± 0.001
ally_health_and_shield			0.15 ± 0.04	0.02 ± 0.005	0.08 ± 0.02	0.002 ± 0.002
ally_health			0.15 ± 0.03	0.02 ± 0.004	0.08 ± 0.02	0.0018 ± 0.0005
ally_distance			0.13 ± 0.03	0.017 ± 0.003	0.07 ± 0.01	-0.0007 ± 0.0004
ally_shield			0.14 ± 0.02	0.018 ± 0.002	0.07 ± 0.01	0.0 ± 0.001
ally_all			0.15 ± 0.04	0.019 ± 0.005	0.08 ± 0.02	0.002 ± 0.001
enemy_health_and_shield			0.43 ± 0.006	0.056 ± 0.002	0.317 ± 0.004	0.038 ± 0.003
enemy_health			0.426 ± 0.006	0.0553 ± 0.001	0.314 ± 0.006	0.037 ± 0.003
enemy_distance			0.18 ± 0.02	0.024 ± 0.002	0.125 ± 0.009	0.006 ± 0.001
enemy_shield			0.14 ± 0.03	0.018 ± 0.003	0.07 ± 0.01	0.0 ± 0.0008
enemy_all			0.54 ± 0.02	0.07 ± 0.002	0.41 ± 0.01	0.052 ± 0.003
nothing			0.14 ± 0.03	0.018 ± 0.004	0.07 ± 0.02	
6h_vs_8z		everything	7.5 ± 0.09	0.87 ± 0.09	0.12 ± 0.02	0.65 ± 0.09
	ally_all_except_actions		0.43 ± 0.02	0.059 ± 0.008	0.31 ± 0.03	0.007 ± 0.002
	ally_health_and_shield		0.41 ± 0.03	0.056 ± 0.009	0.29 ± 0.03	0.0034 ± 0.0007
	ally_health		0.41 ± 0.03	0.055 ± 0.008	0.28 ± 0.03	0.003 ± 0.002
	ally_distance		0.41 ± 0.03	0.056 ± 0.009	0.29 ± 0.03	0.004 ± 0.001
	ally_shield		0.38 ± 0.03	0.051 ± 0.008	0.26 ± 0.03	-0.001 ± 0.001
	ally_all		0.49 ± 0.01	0.067 ± 0.007	0.36 ± 0.02	0.014 ± 0.002
	enemy_health_and_shield		0.45 ± 0.02	0.06 ± 0.008	0.32 ± 0.02	0.008 ± 0.001
	enemy_health		0.46 ± 0.03	0.063 ± 0.01	0.33 ± 0.03	0.011 ± 0.001
	enemy_distance		0.47 ± 0.03	0.064 ± 0.009	0.33 ± 0.03	0.012 ± 0.0004
	enemy_shield		0.39 ± 0.02	0.053 ± 0.008	0.27 ± 0.03	0.001 ± 0.002
	enemy_all		0.55 ± 0.02	0.075 ± 0.008	0.4 ± 0.02	0.0224 ± 0.001
	nothing		0.38 ± 0.03	0.052 ± 0.009	0.27 ± 0.03	
	3s_vs_5z	everything	8.1 ± 0.04	0.3 ± 0.04	0.037 ± 0.007	0.2 ± 0.03
ally_all_except_actions			0.19 ± 0.02	0.024 ± 0.003	0.1 ± 0.02	-0.0014 ± 0.0008
ally_health_and_shield			0.17 ± 0.02	0.021 ± 0.002	0.081 ± 0.01	-0.004 ± 0.001
ally_health			0.17 ± 0.03	0.021 ± 0.004	0.07 ± 0.01	-0.004 ± 0.002
ally_distance			0.2 ± 0.02	0.024 ± 0.002	0.09 ± 0.02	-0.0 ± 0.004
ally_shield			0.18 ± 0.02	0.022 ± 0.002	0.083 ± 0.007	-0.003 ± 0.003
ally_all			0.3 ± 0.08	0.04 ± 0.01	0.17 ± 0.05	0.012 ± 0.007
enemy_health_and_shield			0.26 ± 0.03	0.032 ± 0.004	0.17 ± 0.03	0.007 ± 0.001
enemy_health			0.24 ± 0.03	0.03 ± 0.005	0.15 ± 0.03	0.005 ± 0.002
enemy_distance			0.19 ± 0.02	0.024 ± 0.003	0.09 ± 0.02	-0.0015 ± 0.0006
enemy_shield			0.24 ± 0.02	0.03 ± 0.004	0.13 ± 0.01	0.0049 ± 0.0002
enemy_all			0.26 ± 0.03	0.033 ± 0.004	0.17 ± 0.02	0.008 ± 0.001
nothing			0.2 ± 0.02	0.025 ± 0.003	0.08 ± 0.01	
MM2		everything	9.4 ± 0.1	0.7 ± 0.1	0.08 ± 0.02	0.47 ± 0.09
	ally_all_except_actions		0.41 ± 0.07	0.04 ± 0.01	0.23 ± 0.05	0.007 ± 0.003
	ally_health_and_shield		0.38 ± 0.03	0.041 ± 0.005	0.22 ± 0.03	0.003 ± 0.002
	ally_health		0.39 ± 0.04	0.042 ± 0.006	0.21 ± 0.04	0.004 ± 0.002
	ally_distance		0.34 ± 0.05	0.036 ± 0.007	0.19 ± 0.04	-0.0015 ± 0.0007
	ally_shield		0.35 ± 0.06	0.038 ± 0.008	0.2 ± 0.04	-0.0001 ± 0.0009
	ally_all		0.6 ± 0.1	0.07 ± 0.02	0.36 ± 0.09	0.028 ± 0.008
	enemy_health_and_shield		0.43 ± 0.06	0.046 ± 0.008	0.29 ± 0.04	0.008 ± 0.001
	enemy_health		0.44 ± 0.05	0.047 ± 0.007	0.29 ± 0.04	0.0094 ± 0.0005
	enemy_distance		0.33 ± 0.05	0.036 ± 0.006	0.19 ± 0.04	-0.0022 ± 0.0009
	enemy_shield		0.33 ± 0.04	0.035 ± 0.006	0.18 ± 0.04	-0.002 ± 0.002
	enemy_all		0.59 ± 0.08	0.06 ± 0.01	0.39 ± 0.05	0.026 ± 0.004
	nothing		0.35 ± 0.05	0.038 ± 0.007	0.2 ± 0.04	

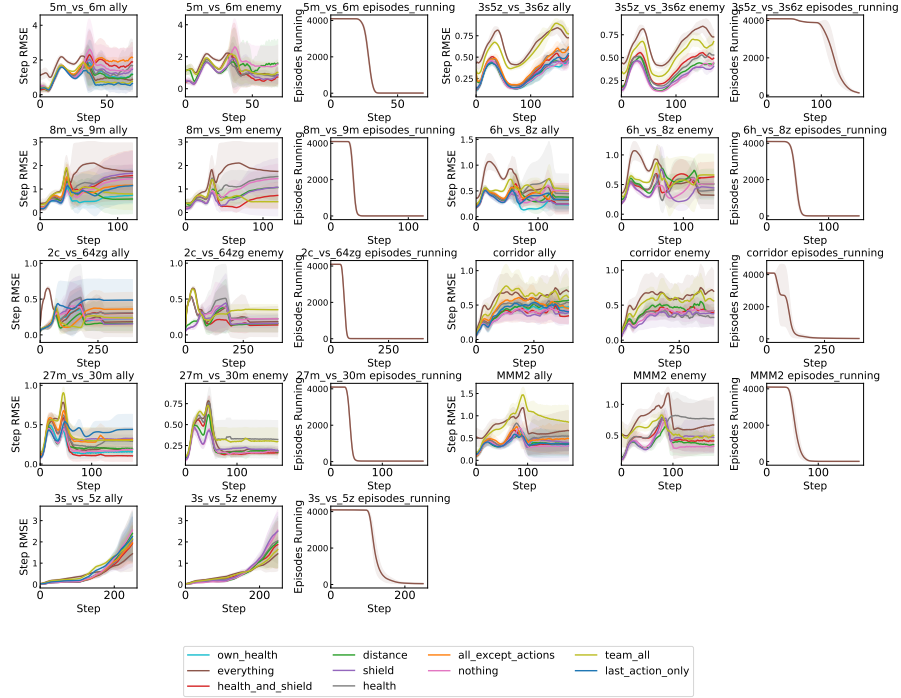


Figure 3: Results of the feature quality experiments for all tested SMAC scenarios not in the main paper.

Table 3: SMACv2 feature quality experiment results

Map	Mask	$\bar{Q}$	$\epsilon_{\text{rmse}}$	$\frac{\epsilon_{\text{rmse}}}{\bar{Q}}$	$\epsilon_{\text{abs}}$	$\frac{\epsilon_{\text{rmse}} - \epsilon_{\text{abs}}}{\bar{Q}}$
5_gen_protoss	everything	$8.4 \pm 0.05$	$1.5 \pm 0.05$	$0.179 \pm 0.006$	$1.18 \pm 0.04$	$0.119 \pm 0.004$
	ally_all_except_actions		$0.69 \pm 0.02$	$0.082 \pm 0.004$	$0.52 \pm 0.01$	$0.022 \pm 0.002$
	ally_health_and_shield		$0.62 \pm 0.02$	$0.074 \pm 0.004$	$0.45 \pm 0.02$	$0.014 \pm 0.002$
	ally_health		$0.54 \pm 0.03$	$0.064 \pm 0.005$	$0.39 \pm 0.02$	$0.004 \pm 0.001$
	ally_distance		$0.52 \pm 0.04$	$0.062 \pm 0.006$	$0.38 \pm 0.03$	$0.0019 \pm 0.0005$
	ally_shield		$0.52 \pm 0.04$	$0.062 \pm 0.006$	$0.38 \pm 0.03$	$0.0018 \pm 0.0008$
	ally_all		$0.85 \pm 0.009$	$0.102 \pm 0.004$	$0.656 \pm 0.01$	$0.042 \pm 0.002$
	enemy_health_and_shield		$0.79 \pm 0.03$	$0.094 \pm 0.004$	$0.6 \pm 0.03$	$0.034 \pm 0.003$
	enemy_health		$0.69 \pm 0.05$	$0.082 \pm 0.007$	$0.51 \pm 0.04$	$0.022 \pm 0.003$
	enemy_distance		$0.56 \pm 0.03$	$0.067 \pm 0.006$	$0.41 \pm 0.03$	$0.0074 \pm 0.0002$
	enemy_shield		$0.57 \pm 0.03$	$0.068 \pm 0.005$	$0.42 \pm 0.02$	$0.008 \pm 0.002$
	enemy_all		$0.96 \pm 0.02$	$0.116 \pm 0.004$	$0.74 \pm 0.02$	$0.055 \pm 0.002$
	nothing			$0.5 \pm 0.03$	$0.06 \pm 0.006$	$0.37 \pm 0.03$
5_gen_terrán	everything	$8.2 \pm 0.1$	$2.0 \pm 0.1$	$0.243 \pm 0.009$	$1.6 \pm 0.1$	$0.148 \pm 0.008$
	ally_all_except_actions		$1.0 \pm 0.05$	$0.1221 \pm 0.001$	$0.75 \pm 0.04$	$0.028 \pm 0.002$
	ally_health_and_shield		$0.82 \pm 0.04$	$0.1 \pm 0.001$	$0.6 \pm 0.03$	$0.0052 \pm 0.0009$
	ally_health		$0.83 \pm 0.04$	$0.1007 \pm 0.0003$	$0.6 \pm 0.04$	$0.0061 \pm 0.0007$
	ally_distance		$0.78 \pm 0.04$	$0.0945 \pm 9e-05$	$0.56 \pm 0.03$	$0.0002 \pm 0.0007$
	ally_shield		$0.75 \pm 0.02$	$0.092 \pm 0.002$	$0.55 \pm 0.02$	$-0.003 \pm 0.001$
	ally_all		$1.18 \pm 0.06$	$0.145 \pm 0.002$	$0.92 \pm 0.06$	$0.05 \pm 0.003$
	enemy_health_and_shield		$0.98 \pm 0.06$	$0.12 \pm 0.003$	$0.72 \pm 0.05$	$0.025 \pm 0.003$
	enemy_health		$0.97 \pm 0.07$	$0.118 \pm 0.003$	$0.71 \pm 0.05$	$0.023 \pm 0.003$
	enemy_distance		$0.82 \pm 0.04$	$0.1002 \pm 0.0002$	$0.59 \pm 0.03$	$0.0055 \pm 0.0006$
	enemy_shield		$0.75 \pm 0.05$	$0.091 \pm 0.005$	$0.54 \pm 0.04$	$-0.004 \pm 0.004$
	enemy_all		$1.2 \pm 0.1$	$0.151 \pm 0.007$	$0.95 \pm 0.09$	$0.057 \pm 0.008$
	nothing			$0.78 \pm 0.03$	$0.0945 \pm 0.0007$	$0.56 \pm 0.03$
5_gen_zerg	everything	$7.2 \pm 0.05$	$2.41 \pm 0.05$	$0.33 \pm 0.02$	$1.86 \pm 0.04$	$0.176 \pm 0.009$
	ally_all_except_actions		$1.36 \pm 0.02$	$0.188 \pm 0.008$	$0.99 \pm 0.02$	$0.031 \pm 0.002$
	ally_health_and_shield		$1.15 \pm 0.03$	$0.16 \pm 0.005$	$0.83 \pm 0.03$	$0.002 \pm 0.001$
	ally_health		$1.17 \pm 0.03$	$0.161 \pm 0.005$	$0.83 \pm 0.02$	$0.0037 \pm 0.0008$
	ally_distance		$1.18 \pm 0.02$	$0.164 \pm 0.007$	$0.85 \pm 0.02$	$0.006 \pm 0.001$
	ally_shield		$1.13 \pm 0.02$	$0.157 \pm 0.006$	$0.81 \pm 0.02$	$-0.0005 \pm 0.0003$
	ally_all		$1.46 \pm 0.02$	$0.202 \pm 0.01$	$1.09 \pm 0.02$	$0.044 \pm 0.002$
	enemy_health_and_shield		$1.31 \pm 0.03$	$0.18 \pm 0.008$	$0.95 \pm 0.03$	$0.024 \pm 0.002$
	enemy_health		$1.3 \pm 0.02$	$0.179 \pm 0.007$	$0.94 \pm 0.02$	$0.022 \pm 0.001$
	enemy_distance		$1.24 \pm 0.01$	$0.173 \pm 0.008$	$0.88 \pm 0.02$	$0.014 \pm 0.002$
	enemy_shield		$1.14 \pm 0.02$	$0.159 \pm 0.006$	$0.82 \pm 0.02$	$0.0005 \pm 0.0008$
	enemy_all		$1.67 \pm 0.01$	$0.23 \pm 0.01$	$1.24 \pm 0.02$	$0.074 \pm 0.008$
	nothing			$1.14 \pm 0.02$	$0.157 \pm 0.007$	$0.81 \pm 0.02$

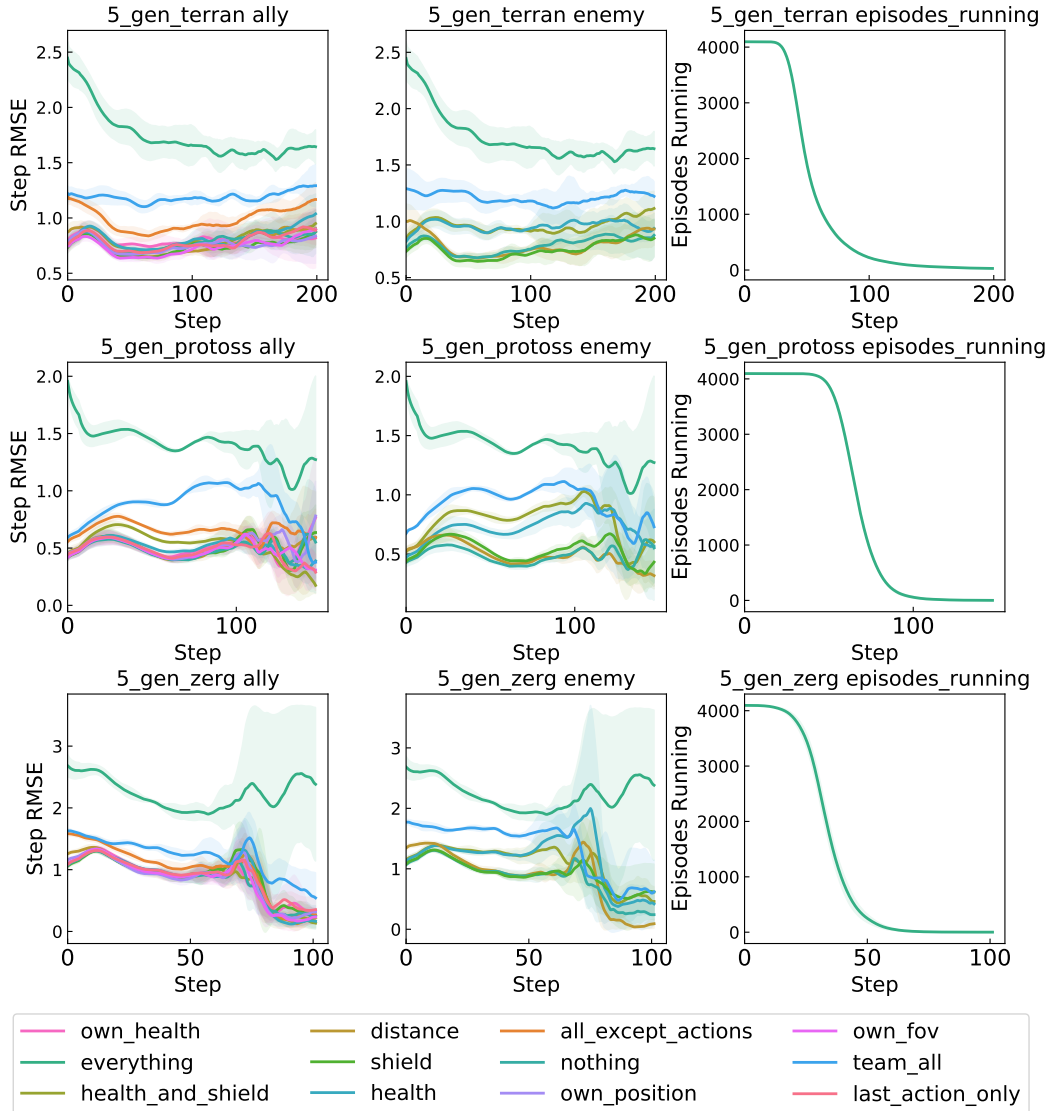


Figure 4: Feature Quality Experiments for the 5 and 10 unit SMACv2 scenarios.

105 of 20 episodes every 2k steps. The hyperparameters used for QMIX are given in Table 4, for MAPPO  
 106 and IPPO in Table 5, and for QPLEX in Table 6.

### 107 3.4 EPO Runs

108 This section provides the implementation details for the Extended Partial Observability challenge in  
 109 SMACv2 and outlines the corresponding training procedure.

110 For the EPO baselines in Section 7.2, we used the same hyperparameters as Section 7.1. Imple-  
 111 mentations of MAPPO and QMIX were also consistent with those in Section 7.1. Each training  
 112 run and corresponding evaluations were conducted across 3 random seeds. Experiments for this  
 113 section were conducted on 80-core CPU machines with NVIDIA GeForce RTX690 2080 Ti or Tesla  
 114 V100-SXM2-16GB GPUs. Both MAPPO and QMIX experiments took between 36-48 hours to  
 115 complete, each running on 8 GPUs, for 3  $p$  values, across 3 seeds and 3 races.

Table 4: QMIX hyperparameters used for experiments. Parameters with (SMAC) or (SMACv2) after them denote that parameter setting was only used for SMAC or SMACv2 experiments respectively. These are the values in the corresponding configuration file in PyMarl[9, 3]. Mac is the code responsible for marshalling inputs to the neural networks, learner is the code used for learning and runner determines whether experience is collected in serial or parallel.

Parameter Name	Value
Action Selector	epsilon greedy
$\epsilon$ Start	1.0
$\epsilon$ Finish	0.05
$\epsilon$ Anneal Time	100000
Runner	parallel
Batch Size Run	4
Buffer Size	5000
Batch Size	128
Optimizer	Adam
$t_{\max}$	10050000
Target Update Interval	200
Mac	n_mac
Agent	n_rnn
Agent Output Type	q
Learner	nq_learner
Mixer	qmix
Mixing Embed Dimension	32
Hypernet Embed Dimension	64
Learning Rate	0.001
$\lambda$ (SMAC)	0.6
$\lambda$ (SMACv2)	0.4

## 116 4 Environment Additional Details

### 117 4.1 SMACv2 Additional Details

118 In this section we describe the generation of teams in SMACv2. For each race, 3 different types  
 119 of units are used. Each race has a special unit that should not be generated too often. For Protoss  
 120 this is the colossus. This is a very powerful unit. If a team has many colossi, the battle will devolve  
 121 into a war about who can use their colossi most effectively. Terran has the medivac unit. This is a  
 122 healing unit that cannot attack the enemy and so is only spawned sparingly. Zerg has the baneling  
 123 unit. This is a suicidal unit which deals area-of-effect damage to enemy units by rolling into them and  
 124 exploding. In scenarios with many banelings, the agents learn to spread out and hide in the corners  
 125 with the hope that the enemy banelings explode and the allies win by default. All of these special  
 126 units are spawned with a probability of 10%. The other units used spawn with a probability of 45%.  
 127 This is summarised in Table 9.

128 There are two changes to the observation space from SMAC. First, each agent observes their own  
 129 field-of-view direction. Secondly, each agent observes their own position in the map as x- and  
 130 y-coordinates. This is normalised by dividing by the map width and height respectively. The only  
 131 change to the state from SMAC was to add the field-of-view direction of each agent to the state.

132 Additionally, we made one small change to the reward function in SMACv2. This was to fix a bug  
 133 where the enemies healing can give allied units reward. There are more details available about this  
 134 problem in the associated Github issue. SMACv2 also has an identical API to the original SMAC,  
 135 allowing for very simple transition between the two frameworks.

136 The code for SMACv2 can be found in the Github repo, where there is a README detailing how to  
 137 run the benchmark with random agents.



Table 5: MAPPO and IPPO hyperparameters used for the experiments on SMAC and SMACv2. Parameters with (SMAC) or (SMACv2) after them denote that parameter setting was only used for SMAC or SMACv2 experiments respectively.

Hyperparameter	Value
Action Selector	multinomial
Mask Before Softmax	True
Runner	parallel
Buffer Size	8
Batch Size Run	8
Batch Size	8
Target Update Interval	200
Learning Rate Actor (SMAC)	0.001
Learning Rate Actor (SMACv2)	0.0005
Learning Rate Critic	0.001
$\tau$	0.995
$\lambda$	0.99
Agent Output Type	pi_logits
Learner	trust_region_learner
Critic (MAPPO)	centralV_rnn_critic
Critic (IPPO)	decentral_rnn_critic
Detach Every	20
Replace Every	None
Mini Epochs Actor	10
Mini Epochs Critic	10
Entropy Loss Coeff	0.0
Advantage Calc Method	GAE
Bootstrap Timeouts	False
Surrogate Clipped	True
Clip Range	0.1
Is Advantage Normalized	True
Observation Normalized	True
Use Popart	False
Normalize Value	False
Obs Agent Id	True
Obs Last Action	False

## 138 4.2 EPO Additional Details

139 In EPO, the first ally agent to observe an enemy is guaranteed to observe the usual SMAC features  
140 associated with that enemy. Upon an enemy being observed for the first time, by any agent on the  
141 ally team, a random binary draw occurs for all other agents (i.e., those that had not yet observed this  
142 particular enemy). Random tie breaking ensures only one agent is guaranteed to see the enemy should  
143 two or more observe it for the first time on the same timestep. The draw is weighted by a tunable  
144 environment parameter  $p$  corresponding to the probability of success. If the draw is successful for a  
145 particular agent, any future instances of the agent observing the enemy will occur as normal, without  
146 masking. If the draw is unsuccessful, that agent will not be able to observe the enemy on future  
147 timesteps, irrespective of whether or not it is within sight range. If the first agent to have observed the  
148 enemy dies, the next time the enemy falls within an ally agent’s sight range that agent is guaranteed  
149 to see the enemy and the random draw occurs again for all other agents.

## 150 5 Analysis of Changes in SMACv2

151 To investigate the impact of the changes introduced in SMACv2, we perform three ablation studies. We  
152 focus only on the 5 unit scenarios, and only train MAPPO on the ablations because it is significantly  
153 faster to train than QMIX. We use the same hyperparameters and setup as in Section 7.1. We ablate  
154 different team compositions by using a single unit type for each of the non-special units in each

Table 6: QPLEX hyperparameters used for experiments on SMACv2.

Hyperparameter	Value
Action Selector	epsilon_greedy
Epsilon Start	1.0
Epsilon Finish	0.05
Epsilon Anneal Time	50000
Runner	parallel
Batch Size Run	4
Buffer Size	5000
Batch Size	128
Optimizer	adam
Target Update Interval	200
Agent Output Type	q
Learner	dmaq_qatten_learner
Double Q	True
Mixer	dmaq_qatten
Mixing Embed Dim	32
Hypernet Embed	64
Adv Hypernet Layers	1
Adv Hypernet Embed	64
Td Lambda	0.8
Num Kernel	4
Is Minus One	True
Is Adv Attention	True
Is Stop Gradient	True
N Head	4
Attend Reg Coef	0.001
State Bias	True
Mask Dead	False
Weighted Head	False
Nonlinear	False
Burn In Period	100
Name	qplex_qatten_sc2

Table 7: QMIX hyperparameter tuning grid

Hyperparameter	Values Tried
$\epsilon$ annealing time	$[100 \times 10^3, 500 \times 10^3]$
$\lambda$	$[0.6, 0.8, 0.4]$

Table 8: MAPPO hyperparameter tuning grid

Hyperparameter	Value
Actor Learning Rate	$[0.0007, 0.0004, 0.0001]$
Clip Range	$[0.15, 0.05, 0.1]$

Table 9: Unit types used per-race in the SMACv2 scenarios.

Race	Unit Types	Probability of Generation
Terran	Marine	0.45
	Marauder	0.45
	Medivac	0.1
Zerg	Zergling	0.45
	Baneling	0.1
	Hydralisk	0.45
Protoss	Stalker	0.45
	Zealot	0.45
	Colossus	0.1

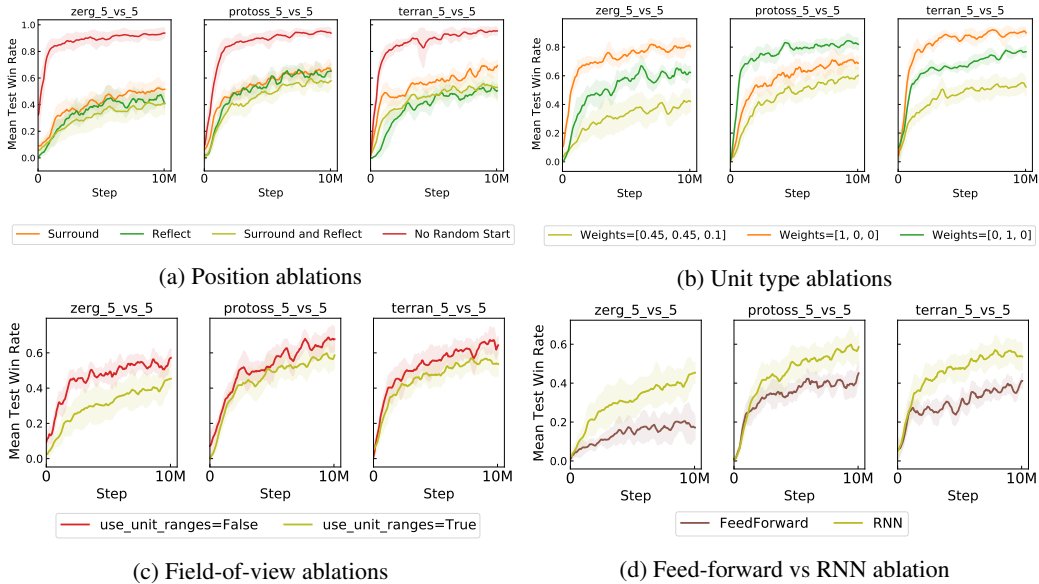


Figure 5: Ablation of SMACv2 unit type, field-of-view and starting position features. Figure 5d shows the performance of a feed-forward policy compared to an RNN baseline. The unit weights are relative to a fixed order of units. For Zerg this is zergling, hydralisk, baneling, for Terran it is marine, marauder, medivac and for Protoss it is stalkers, zealots and colossi. Plots show mean and standard deviation across three seeds.

155 race. We investigate random start positions by only using the *surround* or *reflect* scenarios and by  
 156 removing the random start positions entirely. We also ablate the effect of using the true unit ranges.  
 157 Additionally, we investigate partial observability by comparing feed-forward and recurrent networks.

158 The results of these experiments are shown in Figure 5. The position ablations show that stochastic  
 159 starting positions and hence stochasticity itself contributes greatly to the challenge of SMACv2.  
 160 Without random start positions, MAPPO achieves win rates of close to 100%. Both *surround* and  
 161 *reflect* scenarios have a similar win rate, suggesting similar difficulty. There does not appear to be  
 162 additional difficulty from combining the two. This is logical as the scenarios can be disambiguated  
 163 using start positions at the beginning of the episode.

164 Our experiments of unit type ablations, shown in Figure 5b, indicate that unit variety significantly  
 165 impacts the difficulty of the tasks. All three races show large differences between the easiest unit  
 166 distribution and the baseline, suggesting that a diverse range of units contributes to SMACv2’s  
 167 difficulty. In both the Zerg and Protoss scenarios, the melee-only scenarios are easier than the ablation  
 168 with ranged units. From observing episodes, the enemy AI tends to aggressively pursue allied units,

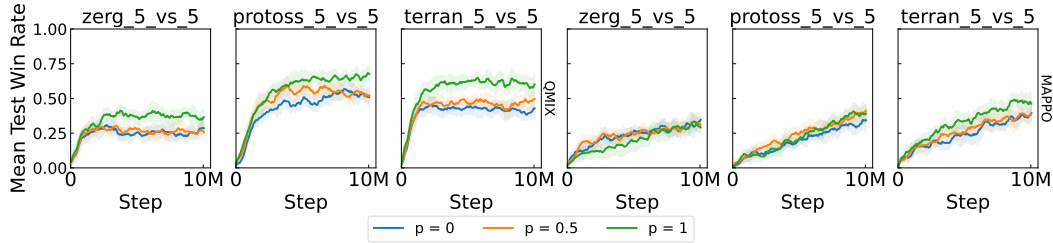


Figure 6: Comparison of the mean test win rate when  $p = 0, 0.5, 1$  for QMIX (left) and MAPPO (right) on SMACv2. Here, the available action mask is still present.

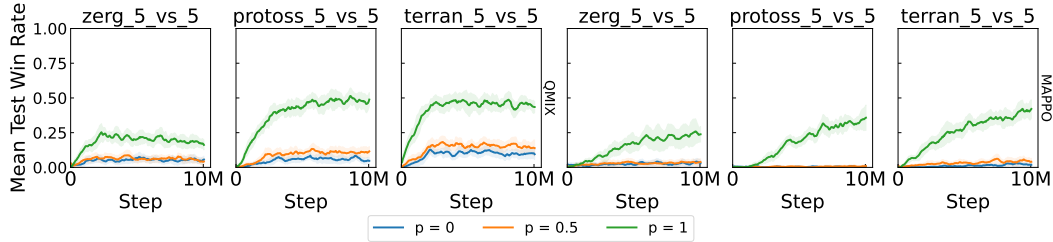


Figure 7: Comparison of mean test win rate when  $p = 0, 0.5, 1$  with 5 agents against 5 enemy units, for QMIX (left) and MAPPO (right). Here, the available actions mask has been removed.

169 which allows the allies to lure it out of position. This is easier to exploit with melee units than ranged  
 170 ones. Terran scenarios have no melee units, but the marine scenarios are slightly easier. Overall these  
 171 results show that generalising to different unit types is a significant part of the challenge in SMACv2.

172 Our field-of-view ablations, shown in Figure 5c, compare the fixed unit sight and attack ranges to  
 173 the SMACv2 versions. There is an increase in difficulty from varying the unit ranges, but this effect  
 174 is small. This suggests that the difficulty of SMACv2 is better explained by the diversity of start  
 175 positions and unit types than the change to the true sight and attack ranges.

176 Finally, to evaluate the effect that the changes had on the partial observability of SMACv2, we  
 177 compare a feedforward network with the performance of the RNN baseline, shown in Figure 5d.  
 178 There is a significant decrease in performance from using the feedforward network. However, the  
 179 feedforward network can still learn reasonable policies. This suggests that although being able to  
 180 resolve partial observability is useful in SMACv2, the primary difficulty consists in generalising to  
 181 a range of micro-management scenarios. This is in contrast to the extended partial observability  
 182 challenge, where resolving partial observability through communication is the primary difficulty. The  
 183 zerg scenario has a larger relative performance difference. Partial observability may be more of a  
 184 problem here because the splash damage done to allies by banelings creates an incentive for allies to  
 185 spread out more.

186 Overall, the results in Figures 5a and 5b show that current MARL methods struggle with the increased  
 187 stochasticity due to map randomisation. To solve this problem, agents must be able to use their  
 188 observations to make inferences about relevant state or joint observation features in the Dec-POMDP  
 189 in more general settings. Current inference capabilities could be improved by research on more  
 190 effective sequence models, such as specialised transformers [14, 6, 1]. The work in this area either  
 191 applies transformer models to traditional RL methods [5], or focuses on paradigms similar to upside-  
 192 down RL [10, 1]. While the latter work shows promising generalisation [8], it is mostly in the offline  
 193 setting. Hu et al. [4] demonstrate a promising sequence model that can be applied to any MARL  
 194 algorithm, but only demonstrate their work on marine units, and rely on the Dec-POMDP having a  
 195 certain structure. Expanding such work to handle more diverse scenarios is an interesting avenue of  
 196 future work.

197 The code for the ablations can be found here. This is distributed under an MIT license. The changes  
 198 on this branch enable easy running of the ablations and add associated environment configurations  
 199 for this purpose.

Table 10: SMACv2 Scenarios with their number of allies, number of enemies, and unit types.

Scenario Name	Number of Allies	Number of Enemies	Unit Types
protoss_5_vs_5	5	5	Stalker, Zealot, Colossus
protoss_10_vs_10	10	10	Stalker, Zealot, Colossus
protoss_20_vs_20	20	20	Stalker, Zealot, Colossus
protoss_10_vs_11	10	11	Stalker, Zealot, Colossus
protoss_20_vs_23	20	23	Stalker, Zealot, Colossus
terran_5_vs_5	5	5	Marine, Marauder, Medivac
terran_10_vs_10	10	10	Marine, Marauder, Medivac
terran_20_vs_20	20	20	Marine, Marauder, Medivac
terran_10_vs_11	10	11	Marine, Marauder, Medivac
terran_20_vs_23	20	23	Marine, Marauder, Medivac
zerg_5_vs_5	5	5	Zergling, Hydralisk, Baneling
zerg_10_vs_10	10	10	Zergling, Hydralisk, Baneling
zerg_20_vs_20	20	20	Zergling, Hydralisk, Baneling
zerg_10_vs_11	10	11	Zergling, Hydralisk, Baneling
zerg_20_vs_23	20	23	Zergling, Hydralisk, Baneling

## 200 6 Limitations and Broader Impact

201 The main limitation of SMACv2 is that it is confined to scenarios within the game of StarCraft II.  
 202 This is an environment which, while complex, cannot represent the dynamics of all multi-agent tasks.  
 203 Evaluation of MARL algorithms therefore should not be limited to one benchmark, but should target  
 204 variety with a range of tasks.

205 Whilst the scenarios in SMACv2 involve battles between two armies of units, only one side can be  
 206 controlled by RL agents. It is technically possible to control two armies using two StarCraft II clients  
 207 that communicate via LAN, which would allow to train two groups of decentralised agents against  
 208 each other, e.g., via self-play. We leave the implementation of this functionality as future work.

209 **SMACv2 aims to contribute to the development of MARL algorithms. As with any machine learning**  
 210 **field, it is possible that improving the capabilities of these algorithms could lead to unethical uses.**  
 211 **However, there are also many potential benefits to better cooperative AI, such as applications in**  
 212 **automated driving among others. We believe that the potential benefits of developing more capable**  
 213 **and cooperative AI outweigh the potential risks.**

## 214 References

- 215 [1] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter  
 216 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning  
 217 via sequence modeling. *Advances in neural information processing systems*, 34, 2021.
- 218 [2] Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS  
 219 Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft  
 220 multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- 221 [3] Jian Hu, Siyang Jiang, Seth Austin Harding, Haibin Wu, and SW Liao. Rethinking the  
 222 implementation tricks and monotonicity constraint in cooperative multi-agent reinforcement  
 223 learning. *arXiv preprint arXiv:2102.03479*, 2021.
- 224 [4] Siyi Hu, Fengda Zhu, Xiaojun Chang, and Xiaodan Liang. Updet: Universal multi-agent rl via  
 225 policy decoupling with transformers. In *International Conference on Learning Representations*,  
 226 2020.
- 227 [5] Shariq Iqbal, Christian A Schroeder De Witt, Bei Peng, Wendelin Boehmer, Shimon Whiteson,  
 228 and Fei Sha. Randomized entity-wise factorization for multi-agent reinforcement learning. In  
 229 Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on*

- 230 *Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 4596–  
231 4606. PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/iqbal21a.](https://proceedings.mlr.press/v139/iqbal21a.html)  
232 [html](https://proceedings.mlr.press/v139/iqbal21a.html).
- 233 [6] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big  
234 sequence modeling problem. *Advances in neural information processing systems*, 34, 2021.
- 235 [7] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
236 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*  
237 *arXiv:1312.5602*, 2013.
- 238 [8] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov,  
239 Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al.  
240 A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- 241 [9] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas  
242 Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon White-  
243 son. The starcraft multi-agent challenge. In *Proceedings of the 18th International Conference*  
244 *on Autonomous Agents and MultiAgent Systems*, pages 2186–2188, 2019.
- 245 [10] Juergen Schmidhuber. Reinforcement learning upside down: Don’t predict rewards—just map  
246 them to actions. *arXiv preprint arXiv:1912.02875*, 2019.
- 247 [11] Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. Qtran:  
248 Learning to factorize with transformation for cooperative multi-agent reinforcement learning.  
249 In *International Conference on Machine Learning*, pages 5887–5896. PMLR, 2019.
- 250 [12] Mingfei Sun, Vitaly Kurin, Guoqing Liu, Sam Devlin, Tao Qin, Katja Hofmann, and Shimon  
251 Whiteson. You may not need ratio clipping in ppo. *arXiv preprint arXiv:2202.00079*, 2022.
- 252 [13] Mingfei Sun, Sam Devlin, Jacob Beck, Katja Hofmann, and Shimon Whiteson. Trust region  
253 bounds for decentralized ppo under non-stationarity. In *Proceedings of the 2023 International*  
254 *Conference on Autonomous Agents and Multiagent Systems*, AAMAS ’23, page 5–13, Richland,  
255 SC, 2023. International Foundation for Autonomous Agents and Multiagent Systems. ISBN  
256 9781450394321.
- 257 [14] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
258 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*  
259 *processing systems*, 30, 2017.
- 260 [15] Jianhao Wang, Zhizhou Ren, Terry Liu, Yang Yu, and Chongjie Zhang. Qplex: Duplex dueling  
261 multi-agent q-learning. In *International Conference on Learning Representations*, 2020.
- 262 [16] Chao Yu, Akash Velu, Eugene Vinitisky, Yu Wang, Alexandre Bayen, and Yi Wu. The surprising  
263 effectiveness of ppo in cooperative, multi-agent games. *arXiv preprint arXiv:2103.01955*, 2021.