



Motivation

- Deep multi-agent reinforcement learning for cooperative multi-agent tasks such as autonomous driving requires learning of policies that allow agents to coordinate their actions.
- Intractably large joint action spaces, partial observability and communication constraints necessitate the learning of decentralised policies even if learning is centralized.
- The absence of communication channels forces agents to coordinate solely based on *implicit communication* through observations.
- Common knowledge allows agents to coordinate *simultaneously* without the need for explicit communication.

The Challenge: How to construct a deep multi-agent reinforcement learning algorithm that benefits from centralized training in the presence of common knowledge, while allowing for fully decentralized execution?

Common Knowledge

- Common knowledge of a group of agents refers to facts that all members 'know', and that "each individual knows that all other individuals know it, each individual knows that all other individuals know that all the individuals know it, and so on" [Osborne & Rubinstein (1994)]
- Notions for imperfect common knowledge, such as *probabilistic common knowledge* exist [Krasucki 1991].

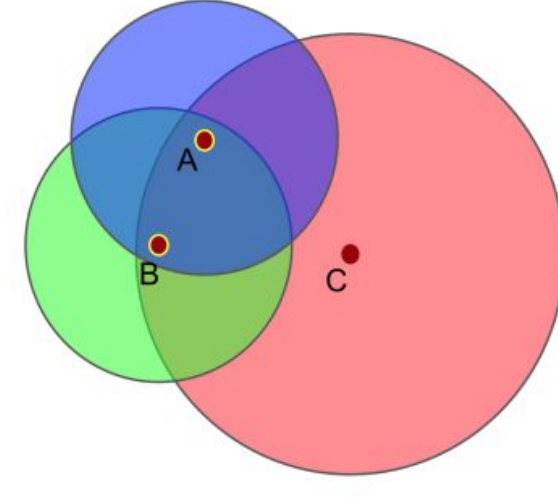


Figure 1: Three agents and their fields of view. A and B's locations are common knowledge to A and B as they are within each other's fields of view. Although C can see A and B, it shares no common knowledge with them.

MACKRL

- Novel actor-critic algorithm with centralised critic and hierarchically factorised joint agent policy.
- Allows end-to-end training while exploiting common knowledge

Algorithm 2 Compute joint policies for a given $\mathbf{u}_{\text{env}}^G \in \mathcal{U}_{\text{env}}^G$ of a group of agents \mathcal{G} in MACKRL

```

function JOINT_POLICY( $\mathbf{u}_{\text{env}}^G | \mathcal{G}, \{\tau_t^a\}_{a \in \mathcal{G}}, \xi$ )
     $a' \sim \mathcal{G}$ ;  $\mathbf{I}^G := \mathcal{I}^G(\tau_t^{a'}, \xi)$   $\triangleright$  random seed in  $\xi$  is common knowledge
     $p_{\text{env}} := 0$   $\triangleright$  common knowledge  $\mathbf{I}^G$  is identical for every agent  $a' \in \mathcal{G}$ 
     $\triangleright$  initialise probability for choosing environmental joint action  $\mathbf{u}_{\text{env}}^G$ 
    for  $\mathbf{u}^G \in \mathcal{U}_{\text{env}}^G$  do  $\triangleright$  add probability to choose  $\mathbf{u}_{\text{env}}^G$  for all outcomes of  $\pi^G$ 
        if  $\mathbf{u}^G = \mathbf{u}_{\text{env}}^G$  then  $\triangleright$  if  $\mathbf{u}^G$  is the environmental joint action in question
             $p_{\text{env}} := p_{\text{env}} + \pi^G(\mathbf{u}_{\text{env}}^G | \mathbf{I}^G)$ 
        if  $\mathbf{u}^G \neq \mathbf{u}_{\text{env}}^G$  then  $\triangleright$  if  $\mathbf{u}^G = \{\mathcal{G}^1, \dots, \mathcal{G}^k\}$  is a set of disjoint subgroups
             $p_{\text{env}} := p_{\text{env}} + \pi^G(\mathbf{u}^G | \mathbf{I}^G) \prod_{\mathcal{G}' \in \mathcal{G}^k} \text{JOINT\_POLICY}(\mathbf{u}_{\text{env}}^{G'} | \mathcal{G}', \{\tau_t^a\}_{a \in \mathcal{G}'}, \xi)$ 
    return  $p_{\text{env}}$   $\triangleright$  return probability that controller  $\pi^G$  would have chosen  $\mathbf{u}_{\text{env}}^G$ 

```

Algorithm 1 Decentralised action selection for agent $a \in \mathcal{A}$ in MACKRL

```

function SELECT_ACTION( $a, \tau_t^a, \xi$ )
     $\mathcal{G} := \mathcal{A}$   $\triangleright$  random seed in  $\xi$  is common knowledge
     $\triangleright$  initialise the group  $\mathcal{G}$  of all agents
     $\mathbf{u}_t^G \sim \pi^G(\cdot | \mathcal{I}^G(\tau_t^a, \xi))$   $\triangleright$   $\mathbf{u}_t^G$  is either a joint environmental action in  $\mathcal{U}_{\text{env}}^G$ ...
    while  $\mathbf{u}_t^G \notin \mathcal{U}_{\text{env}}^G$  do  $\triangleright$  ... or a set of disjoint subgroups  $\{\mathcal{G}^1, \dots, \mathcal{G}^k\}$ 
         $\mathcal{G} := \{\mathcal{G}' | a \in \mathcal{G}', \mathcal{G}' \in \mathcal{U}_t^G\}$   $\triangleright$  select subgroup containing agent  $a$ 
         $\mathbf{u}_t^G \sim \pi^G(\cdot | \mathcal{I}^G(\tau_t^a, \xi))$   $\triangleright$  draw an action for that subgroup
    return  $\mathbf{u}_t^G$   $\triangleright$  return environmental action  $\mathbf{u}_t^G \in \mathcal{U}_{\text{env}}^G$  of agent  $a$ 

```

Pairwise MACKRL

Pairwise MACKRL restricts the joint policy to a three-level hierarchical controller, where the top level is a pair selector conditioning on common knowledge between all the agents, the second level consists of pair controllers conditioning on pair-wise common knowledge. The second level can choose to delegate to the lowest level, which consists of independent learners.

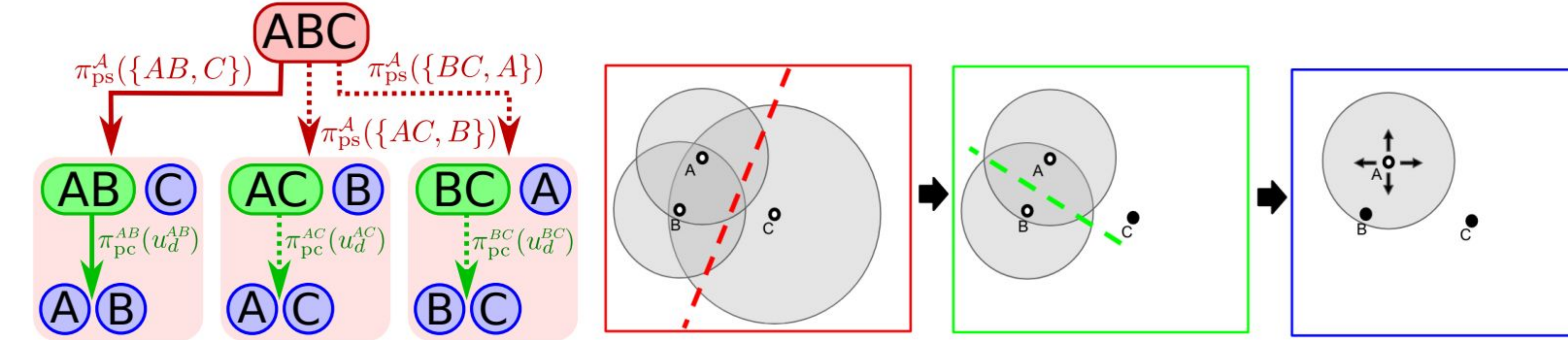


Figure 2: An example for Pairwise MACKRL. [Left]: the full hierarchy for 3 agents (dependencies on common knowledge are omitted for clarity). Only solid arrows are computed during decentralised sampling with Algorithm 1, while all arrows must be computed recursively during centralised training (see Algorithm 2). [Right]: the (maximally) 3 steps of decentralized sampling from the perspective of agent A. (i) Pair selector π_{ps}^A chooses the partition $\{AB, C\}$ based on the common knowledge of all agents $\mathcal{I}^{ABC}(\tau^A, \xi) = \emptyset$. (ii) Based on the common knowledge of pair A and B, $\mathcal{I}^{AB}(\tau^A, \xi)$, the pair controller π_{pc}^{AB} can either choose a joint action $(u_{\text{env}}^A, u_{\text{env}}^B)$, or delegate to individual controllers by selecting u_d^{AB} . (iii) If delegating, the individual controller π^A must select the action u_{env}^A for the single agent A. All steps can be computed based on A's history τ^A .

Level	Policy / Controller	# π
1	$\pi_{\text{ps}}^A(u_{\text{ps}}^A \mathcal{I}_{s_t}^A, u_{t-1}^A, h_{t-1}^A)$	1
2	$\pi_{\text{pc}}^{aa'}(u_{\text{pc}}^{aa'} \mathcal{I}_{s_t}^{aa'}, u_{t-1}^{aa'}, h_{t-1}^{aa'}, aa')$	3
3	$\pi^a(u^a z_t^a, h_{t-1}^a, u_{t-1}^a, a)$	3

Top: Policy hierarchy for Pairwise MACKRL. Right: Explicitly Written-out joint policy for Pairwise MACKRL with three agents.

$$\begin{aligned}
 \pi_{\theta}(u_{\text{env}}^1, u_{\text{env}}^2, u_{\text{env}}^3) = & \pi_{\text{ps}, \theta}^A(u_{\text{ps}}^A = \{1, \{2, 3\}\} | \mathcal{I}_{s_t}^{1,2,3}) \cdot \pi_{\theta}^1(u_{\text{env}}^1 | \tau^1) \\
 & \cdot (\pi_{\text{pc}, \theta}^{2,3}(u_{\text{pc}}^{2,3} | \mathcal{I}_{s_t}^{2,3}) + \pi_{\text{pc}, \theta}^2(u_d^{2,3} | \mathcal{I}_{s_t}^{2,3}) \cdot \pi_{\theta}^2(u_{\text{env}}^2 | \tau^2) \cdot \pi_{\theta}^3(u_{\text{env}}^3 | \tau^3)) \\
 & + \pi_{\text{ps}, \theta}^A(u_{\text{ps}}^A = \{2, \{1, 3\}\} | \mathcal{I}_{s_t}^{1,2,3}) \cdot \pi_{\theta}^2(u_{\text{env}}^2 | \tau^2) \\
 & \cdot (\pi_{\text{pc}, \theta}^{1,3}(u_{\text{pc}}^{1,3} | \mathcal{I}_{s_t}^{1,3}) + \pi_{\text{pc}, \theta}^1(u_d^{1,3} | \mathcal{I}_{s_t}^{1,3}) \cdot \pi_{\theta}^1(u_{\text{env}}^1 | \tau^1) \cdot \pi_{\theta}^3(u_{\text{env}}^3 | \tau^3)) \\
 & + \pi_{\text{ps}, \theta}^A(u_{\text{ps}}^A = \{3, \{1, 2\}\} | \mathcal{I}_{s_t}^{1,2,3}) \cdot \pi_{\theta}^3(u_{\text{env}}^3 | \tau^3) \\
 & \cdot (\pi_{\text{pc}, \theta}^{1,2}(u_{\text{pc}}^{1,2} | \mathcal{I}_{s_t}^{1,2}) + \pi_{\text{pc}, \theta}^{1,2}(u_d^{1,2} | \mathcal{I}_{s_t}^{1,2}) \cdot \pi_{\theta}^1(u_{\text{env}}^1 | \tau^1) \cdot \pi_{\theta}^2(u_{\text{env}}^2 | \tau^2)).
 \end{aligned}$$

Results - Matrix game with common knowledge

- We investigate a two-agent single-step matrix game where the reward matrix is chosen between A and B at random. The choice of matrix played is available as common knowledge with probability p
- We demonstrate that Pairwise MACKRL outperforms both Joint Action-Learning (JAL) and Independent Learning (IL) under decentralised execution if common knowledge is available.

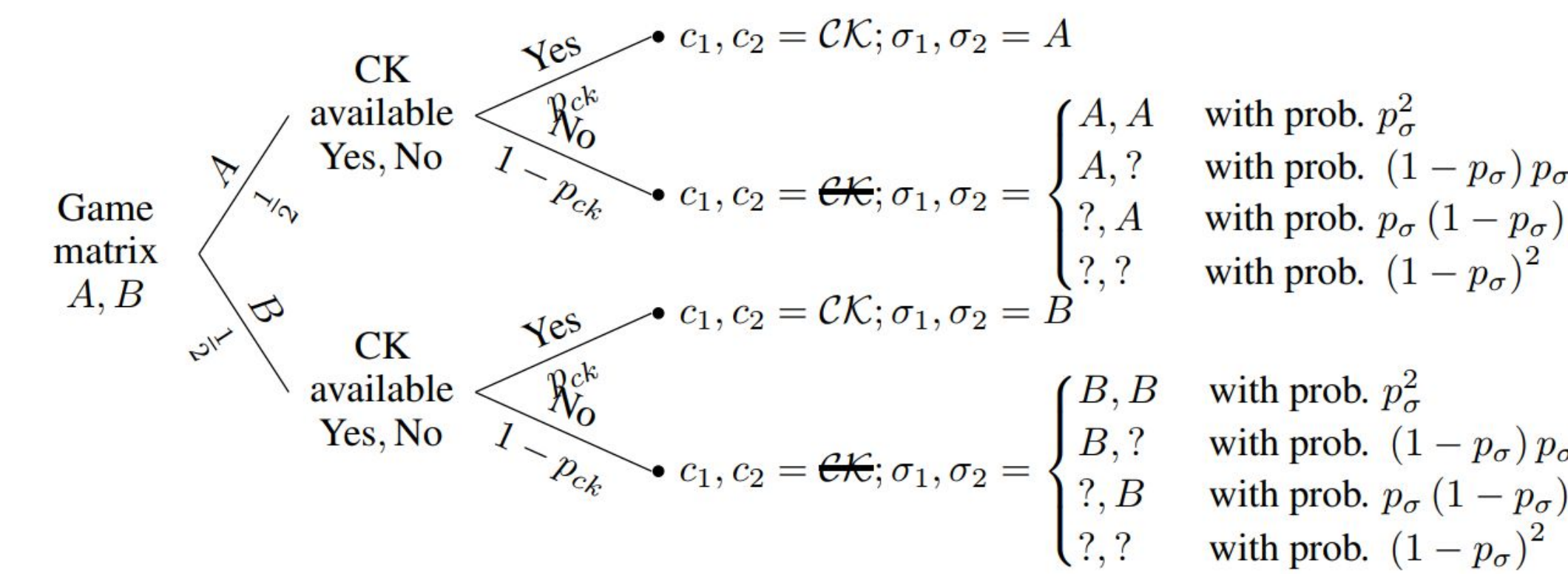
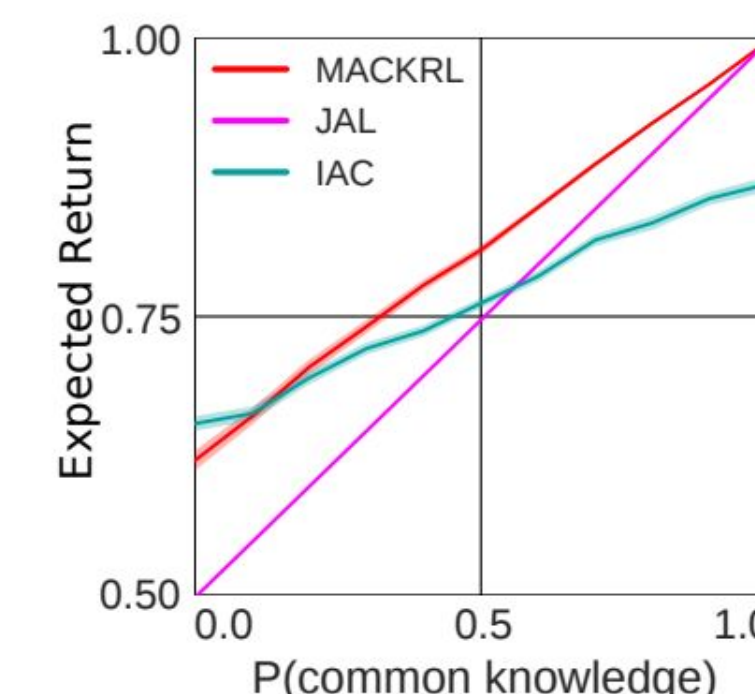


Figure 7: Probability tree for our simple single-step matrix game. The game chooses randomly between matrix A or B, and whether common knowledge is available or not. If common knowledge is available, both agents can condition their actions on the game matrix choice. Otherwise, both agents independently only have a random chance of observing the game matrix choice. Here, p_{ck} is the probability that common knowledge exists and p_o is the probability that an agent independently observes the game matrix choice. The observations of each agent 1 and 2 are given by tuples (c_1, σ_1) and (c_2, σ_2) , respectively, where $c_1, c_2 \in \{\text{CK}, \text{?}\}$ and $\sigma_1, \sigma_2 \in \{A, B, \text{?}\}$.



$$\frac{1}{5} \begin{bmatrix} 0 & 0 & 1 & 0 & 5 \\ 0 & 0 & 2 & 0 & 0 \\ 1 & 2 & 4 & 2 & 1 \\ 0 & 0 & 2 & 0 & 0 \\ 5 & 0 & 1 & 0 & 0 \end{bmatrix} \quad \frac{1}{5} \begin{bmatrix} 5 & 0 & 0 & 2 & 0 \\ 0 & 1 & 2 & 4 & 2 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 5 \end{bmatrix}$$

Above: The two game reward matrices A and B. Left: Expected return of MACKRL, JAL and IAC

Results - StarCraft II

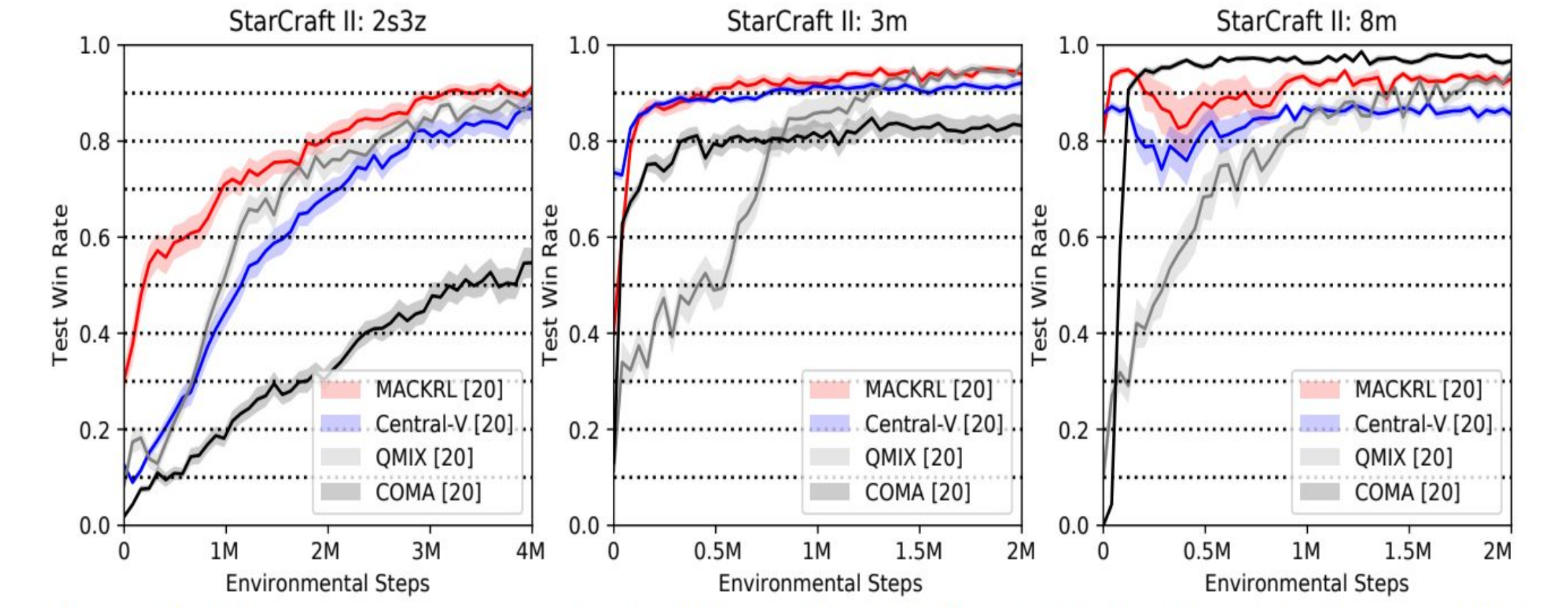


Figure 4: Win rate at test time across StarCraft II scenarios: 2 Stalkers & 3 Zealots [left], 3 Marines [middle] and 8 Marines [right]. Plots show means and their standard errors with [number of runs].

- Challenging unit micromanagement tasks under partial observability constraints with circular field of view.
- MACKRL outperforms its Central-V baseline in sample efficiency and limit performance on all maps.
- MACKRL also significantly outperforms state-of-the-art COMA and QMIX algorithms in terms of sample efficiency.

Scalability, Robustness & Introspection

- MACKRL can exploit *probabilistic* common knowledge arising from noisy sensor using correlated sampling.
- MACKRL can scale to many agents by restricting pair selector to random subsets
- Pair controller introspection verifies use of common knowledge

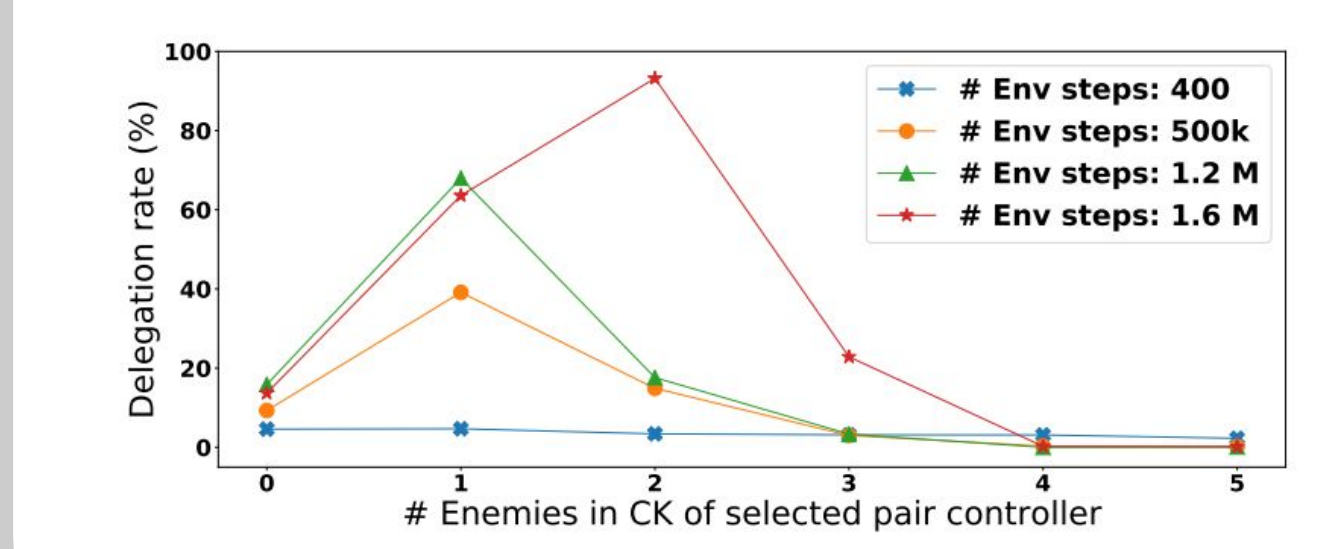
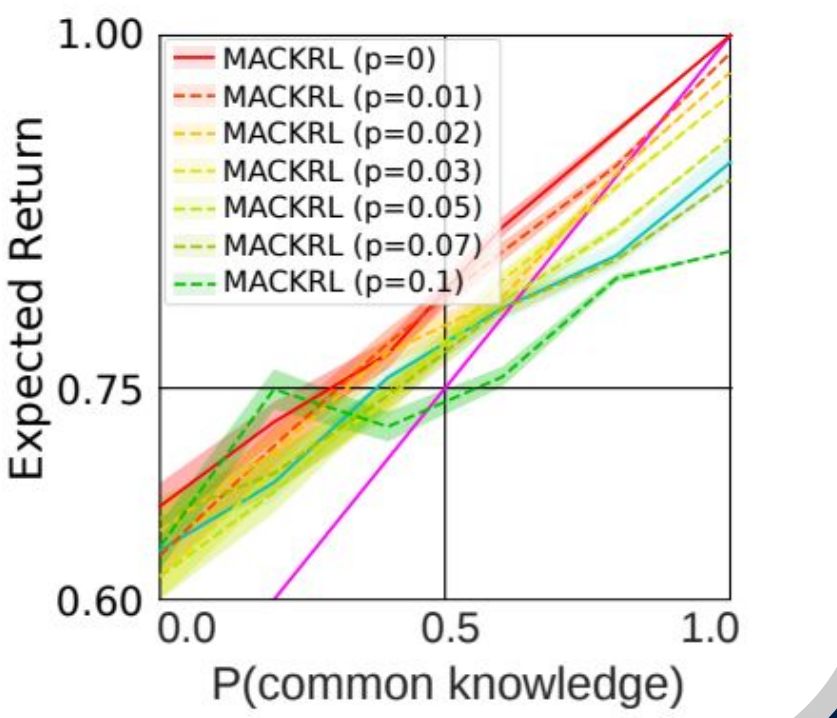


Figure 8: Delegation rate vs. number of enemies (2s3z) in the common knowledge of the pair controller over training.



Discussion and Future Work

- MACKRL uses a novel hierarchical policy structure that allows to exploit common knowledge between agents during fully decentralized execution. This results in outperformance of strong baselines and even state-of-the-art on StarCraft II. We demonstrate MACKRL's ability to scale to many agents and its robustness to sensor noise.
- Future work will address off-policy variants of MACKRL and investigate how limited-bandwidth communication can be exploited in the presence of common knowledge.