

1 A Appendix

2 A.1 Game Environment

3 Figure 1 shows the UI interface of *Honor of Kings*. For fair comparisons, all experiments in this paper
4 were carried out using a fixed released gamecore version (Version 3.73 series) of *Honor of Kings*.



Figure 1: **The UI interface of *Honor of Kings***. The hero controlled by the player is called *Main Hero*. The player controls the hero’s movement through the bottom-left wheel (C.1) and releases the hero’s skills through the bottom-right buttons (C.2, C.3). The player can observe the local view via the screen, observe the global view via the top-left mini-map (A), and obtain game states via the top-right dashboard (B).

5 A.2 In-game Signaling System

6 Figure 2 demonstrates the in-game signaling system of *Honor of Kings*. Players can communicate
7 and collaborate with teammates through the in-game signaling system. In the **Human-AI Game**
8 **Test**, humans can send macro-strategies to agents through signals like A in figure 2, and these signals
9 are displayed to teammates in the form of D. The MCC framework converts these explicit messages,
10 i.e., signals, into meta-commands by the hand-crafted command converter function f^{cc} and broadcast
11 them to all agent teammates. And the MCC framework can also convert the meta-commands sent
12 from agents into signals by the inverse of f^{cc} and broadcast them to all human teammates.

13 Voice (B.2) and text (B.1 and B.3) are two other forms of communication. In the future, we consider
14 introducing a general meta-command encoding model that can handle all forms of explicit messages
15 (signals, voice, and text).



Figure 2: **The in-game signaling system of *Honor of Kings***. Players can send their macro-strategies by dragging signal buttons (A.2) to the corresponding locations (A.1) in the mini-map. The sent result is displayed in the form of a yellow circle (D). C is the convenience signals representing attack, retreat, and assembly, respectively. Voice (B.2) and text (B.1 and B.3) are two other forms of communication.

16 A.3 Hero Pool

17 Table 1 shows the full hero pool and 20 hero pool used in the **Experiments**. Each match involves two
 18 lineups playing against each other, and each lineup consists of five randomly picked heroes.

Table 1: Hero pool used in the **Experiments**.

Full Hero pool	Lian Po, Xiao Qiao, Zhao Yun, Mo Zi, Da Ji, Ying Zheng, Sun Shangxiang, Luban Qihao, Zhuang Zhou, Liu Chan Gao Jianli, A Ke, Zhong Wuyan, Sun Bin, Bian Que, Bai Qi, Mi Yue, Lv Bu, Zhou Yu, Yuan Ge Xia Houdun, Zhen Ji, Cao Cao, Dian Wei, Gongben Wucang, Li Bai, Make Boluo, Di Renjie, Da Mo, Xiang Yu Wu Zetian, Si Mayi, Lao Fuzi, Guan Yu, Diao Chan, An Qila, Cheng Yaojin, Lu Na, Jiang Ziya, Liu Bang Han Xin, Wang Zhaojun, Lan Lingwang, Hua Mulan, Ai Lin, Zhang Liang, Buzhi Huowu, Nake Lulu, Ju Youjing, Ya Se Sun Wukong, Niu Mo, Hou Yi, Liu Bei, Zhang Fei, Li Yuanfang, Yu Ji, Zhong Kui, Yang Yuhuan, Chengji Sihan Yang Jian, Nv Wa, Ne Zha, Ganjiang Moye, Ya Dianna, Cai Wenji, Taiyi Zhenren, Donghuang Taiyi, Gui Guzi, Zhu Geliang Da Qiao, Huang Zhong, Kai, Su Lie, Baili Xuance, Baili Shouyue, Yi Xing, Meng Qi, Gong Sunli, Shen Mengxi Ming Shiyin, Pei Qinhu, Kuang Tie, Mi Laidi, Yao, Yun Zhongjun, Li Xin, Jia Luo, Dun Shan, Sun Ce Zhu Bajie, Shangguan Waner, Ma Chao, Dong Fangyao, Xi Shi, Meng Ya, Luban Dashi, Pan Gu, Chang E, Meng Tian Jing, A Guduo, Xia Luote, Lan, Sikong Zhen, Erin, Yun ying, Jin Chan, Fei, Sang Qi
20 Hero Pool	Jing, Pan Gu, Zhao Yun, Ju Youjing, Donghuang Taiyi, Zhang Fei, Gui Guzi, Da Qiao, Sun Shangxiang, Luban Qihao, Chengji Sihan, Huang Zhong, Zhuang Zhou, Lian Po, Liu Bang, Zhong Wuyan, Yi Xing, Zhou Yu, Xi Shi, Zhang Liang

19 A.4 Agent Action

20 Table 2 shows the action space of agents.

Table 2: The action space of agents.

Action	Detail	Description
What	Illegal action	Placeholder.
	None action	Executing nothing or stopping continuous action.
	Move	Moving to a certain direction determined by move x and move y.
	Normal Attack	Executing normal attack to an enemy unit.
	Skill1	Executing the first skill.
	Skill2	Executing the second skill.
	Skill3	Executing the third skill.
	Skill4	Executing the fourth skill (only a few heroes have Skill4).
	Summoner ability	An additional skill choosing before the game begins (10 to choose).
	Return home(Recall)	Returning to spring, should be continuously executed.
	Item skill	Some items can enable an additional skill to player's hero.
	Restore	Blood recovering continuously in 10s, can be disturbed.
Collaborative skill	Skill given by special ally heroes.	
How	Move X	The x-axis offset of moving direction.
	Move Y	The y-axis offset of moving direction.
	Skill X	The x-axis offset of a skill.
	Skill Y	The y-axis offset of a skill.
Who	Target unit	The game unit(s) chosen to attack.

21 A.5 Reward Design

22 Table 3 demonstrates the details of the designed environment reward.

23 A.6 Infrastructure Design

24 Figure 3 shows the infrastructure of the training system, which consists of four pivotal components:
 25 AI Server, Inference Server, RL Learner, and Memory Pool. The AI Server (the actor) covers the
 26 interaction logic between the agents and the environment. The Inference Server is used for the
 27 centralized batch inference on the GPU side. The RL Learner (the learner) is a distributed training
 28 environment for RL models. And the Memory Pool is used for storing the experience, implemented
 29 as a memory-efficient circular queue.

30 As is known to all, training complex game AI systems often require a large amount of computing
 31 resources, such as AlphaGo Lee Sedol (280 GPUs), OpenAI Five Final (1920 GPUs), and AlphaStar
 32 Final (3072 TPUv3 cores), we also use hundreds of GPUs for training the agents. Another future
 33 work is to improve resource utilization using fewer computing resources.

Table 3: The details of the environment reward.

Head	Reward Item	Weight	Type	Description
Farming Related	Gold	0.005	Dense	The gold gained.
	Experience	0.001	Dense	The experience gained.
	Mana	0.05	Dense	The rate of mana (to the fourth power).
	No-op	-0.00001	Dense	Stop and do nothing.
	Attack monster	0.1	Sparse	Attack monster.
KDA Related	Kill	1	Sparse	Kill a enemy hero.
	Death	-1	Sparse	Being killed.
	Assist	1	Sparse	Assists.
	Tyrant buff	1	Sparse	Get buff of killing tyrant, dark tyrant, storm tyrant.
	Overlord buff	1.5	Sparse	Get buff of killing the overlord.
	Expose invisible enemy	0.3	Sparse	Get visions of enemy heroes.
	Last hit	0.2	Sparse	Last hitting an enemy minion.
Damage Related	Health point	3	Dense	The health point of the hero (to the fourth power).
	Hurt to hero	0.3	Sparse	Attack enemy heroes.
Pushing Related	Attack turrets	1	Sparse	Attack turrets.
	Attack crystal	1	Sparse	Attack enemy home base.
Win/Lose Related	Destroy home base	2.5	Sparse	Destroy enemy home base.

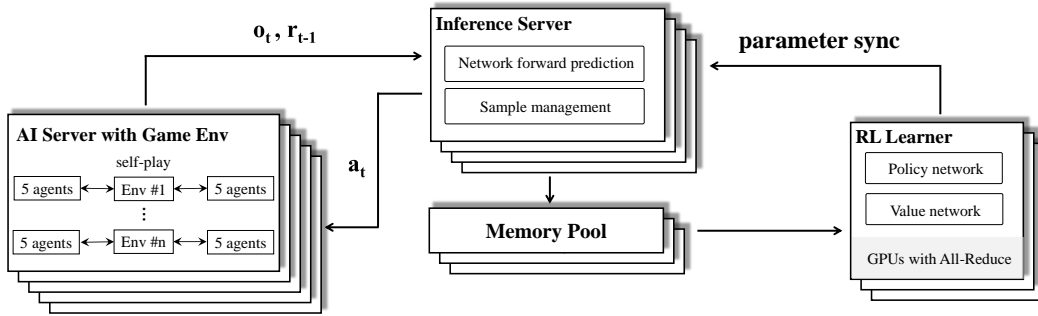


Figure 3: The designed infrastructure.

34 A.7 Feature Design

35 A.7.1 CEN

36 See Table 4.

37 A.7.2 MCCAN

38 See Table 5.

39 A.7.3 Feature of CS

40 See Table 6.

41 A.8 Network Architecture

42 A.8.1 CEN

43 Figure 4 shows the detailed model structure of CEN. The CEN predicts a meta-command Softmax
 44 distribution for each agent based on its current observation. The outputted meta-command indicates
 45 the macro-strategy for future T^{mc} steps.

46 A.8.2 MCCAN

47 Figure 5 shows the detailed model structure of MCCAN. The MCCAN predicts a sequence of actions
 48 for each agent based on its observation and the meta-command sampled from the Top- k Softmax
 49 distribution of CEN. The observations are processed by a deep LSTM, which maintains memory

Table 4: Feature details of CEN.

Feature Class	Field	Description	Dimension
1. Unit feature	Scalar	Includes heroes, minions, monsters, and turrets	3946
Heroes	Status	Current HP, mana, speed, level, gold, KDA, and magical attack and defense, etc.	1562
	Position	Current 2D coordinates	20
Minions	Status	Current HP, speed, visibility, killing income, etc.	920
	Position	Current 2D coordinates	80
Monsters	Status	Current HP, speed, visibility, killing income, etc.	728
	Position	Current 2D coordinates	56
Turrets	Status	Current HP, locked targets, attack speed, etc.	540
	Position	Current 2D coordinates	40
2. In-game stats feature	Scalar	Real-time statistics of the game	104
Static statistics	Time	Current game time	57
	Camp	Types of two camps	1
	Alive heroes	Number of alive heroes of two camps	10
	Kill	Kill number of each camp	6
	Alive turrets	Number of alive turrets of two camps	8
Comparative statistics	Alive heroes diff	Alive heroes difference between two camps	11
	Kill diff	Kill difference between two camps	5
	Alive turrets diff	Alive turrets difference between two camps	6

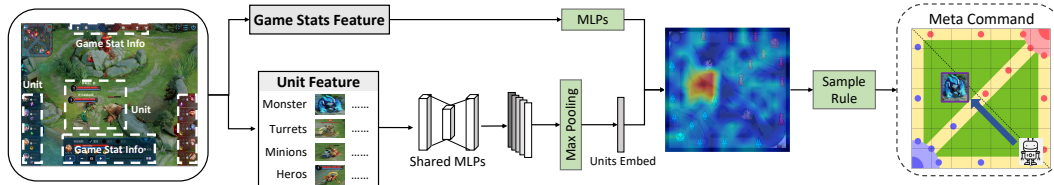


Figure 4: The CEN model structure.

50 among steps. we use the target attention mechanism to improve the accuracy of the model prediction,
 51 and we design the action mask module to eliminate unnecessary actions for efficient exploration.
 52 Additionally, we introduce a value mixer module [4] to model team value for improving the accuracy
 53 of the value estimation. Finally, following [5] and [1], we adopt hierarchical heads of actions,
 54 including three parts: 1) What action to take; 2) who to target; 3) how to act.

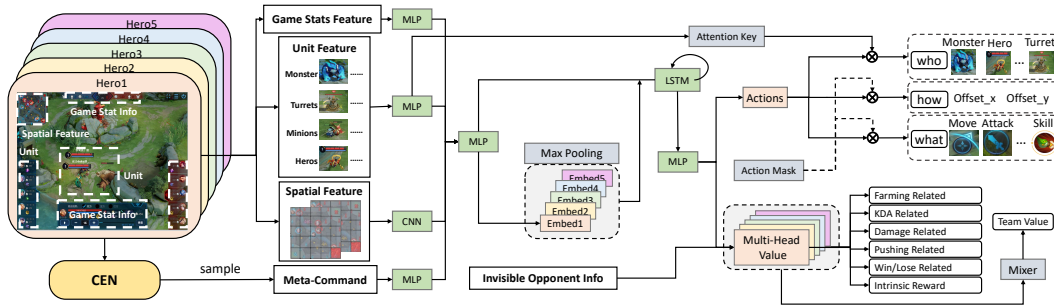


Figure 5: The MCCAN model structure.

55 **A.9 Details of Human-AI Game Test**

56 **A.9.1 Participant**

57 We contacted the game provider and got a test authorization. The game provider found us participants
 58 who meet the requirements. During the **Human-AI Game Test**, we only know the rank-level and
 59 game experience information of participants and do not know their identity information. And special
 60 equipment and game accounts are provided to the participants to prevent the leakage of equipment

Table 5: Feature details of MCCAN.

Feature Class	Field	Description	Dimension
1. Unit feature	Scalar	Includes heroes, minions, monsters, and turrets	8599
Heroes	Status	Current HP, mana, speed, level, gold, KDA, buff, bad states, orientation, visibility, etc.	1842
	Position	Current 2D coordinates	20
	Attribute	Is main hero or not, hero ID, camp (team), job, physical attack and defense, magical attack and defense, etc.	1330
	Skills	Skill 1 to Skill N's cool down time, usability, level, range, buff effects, bad effects, etc.	2095
	Item	Current item lists	60
Minions	Status	Current HP, speed, visibility, killing income, etc.	1160
	Position	Current 2D coordinates	80
	Attribute	Camp (team)	80
	Type	Type of minions (melee creep, ranged creep, siege creep, super creep, etc.)	200
Monsters	Status	Current HP, speed, visibility, killing income, etc.	868
	Position	Current 2D coordinates	56
	Type	Type of monsters (normal, blue, red, tyrant, overlord, etc.)	168
Turrets	Status	Current HP, locked targets, attack speed, etc.	520
	Position	Current 2D coordinates	40
	Type	Type of turrets (tower, high tower, crystal, etc.)	80
2. In-game stats feature	Scalar	Real-time statistics of the game	68
Static statistics	Time	Current game time	5
	Gold	Gold of two camps	12
	Alive heroes	Number of alive heroes of two camps	10
	Kill	Kill number of each camp	6
	Alive turrets	Number of alive turrets of two camps	8
Comparative statistics	Gold diff	Gold difference between two camps	5
	Alive heroes diff	Alive heroes difference between two camps	11
	Kill diff	Kill difference between two camps	5
	Alive turrets diff	Alive turrets difference between two camps	6
3. Invisible opponent information	Scalar	Invisible information used for the value net	560
Opponent heroes	Position	Current 2D coordinates, distances, etc.	120
NPC	Position	Current 2D coordinates of all non-player characters, including minions, monsters, and turrets	440
4. Spatial feature	Spatial	2D image-like, extracted in channels for convolution	6x17x17
Skills	Region	Potential damage regions of ally and enemy skills	2x17x17
	Bullet	Bullets of ally and enemy skills	2x17x17
Obstacles	Region	Forbidden region for heroes to move	1x17x17
Bushes	Region	Bush region for heroes to hide	1x17x17
5. Meta-Command feature	Spatial	Flattened Meta-Command	144

61 and account information. The game statistics we collect are for experimental purposes only and are
62 not disclosed to the public.

63 The participants consisted of 15 strong humans (top 1%) and 15 average humans (top 30%). All
64 participants have more than three years of experience in *Honor of Kings* and promise to be familiar
65 with all mechanics in the game, including the in-game signaling system in Figure 2. We used m AI +
66 n Human mode to evaluate the performance of agents teaming up with different numbers of humans,
67 where $m + n = 5$. Each participant is asked to randomly team up with three different types of agents,
68 including the MC-Base agents, the MC-Rand agents, and the MCC agents. For fair comparisons, we
69 adopt the MC-Base agent as the opponent for all tests. Each participant tested 20 matches for the 4
70 AI + 1 Human mode. Each strong human participant tested additional 10 matches for the 3 AI + 2
71 Human and 2 AI + 3 Human modes, respectively. In all tests, participants were not told the type of
72 agent teammates.

73 In addition, as mentioned in [5, 1], the response time of agents is usually set to 193ms, including
74 observation delay (133ms) and response delay (60ms). The average APMs of agents and top e-sport
75 players are usually comparable (80.5 and 80.3, respectively). To make our test results more accurate,

Table 6: Feature details of CS.

Feature Class	Field	Description	Dimension
1. Unit feature	Scalar	Includes heroes, minions, monsters, and turrets	3946
Heroes	Status	Current HP, mana, speed, level, gold, KDA, and magical attack and defense, etc.	1562
	Position	Current 2D coordinates	20
Minions	Status	Current HP, speed, visibility, killing income, etc.	920
	Position	Current 2D coordinates	80
Monsters	Status	Current HP, speed, visibility, killing income, etc.	728
	Position	Current 2D coordinates	56
Turrets	Status	Current HP, locked targets, attack speed, etc.	540
	Position	Current 2D coordinates	40
2. In-game stats feature	Scalar	Real-time statistics of the game	104
Static statistics	Time	Current game time	57
	Camp	Types of two camps	1
	Alive heroes	Number of alive heroes of two camps	10
	Kill	Kill number of each camp	6
	Alive turrets	Number of alive turrets of two camps	8
Comparative statistics	Alive heroes diff	Alive heroes difference between two camps	11
	Kill diff	Kill difference between two camps	5
	Alive turrets diff	Alive turrets difference between two camps	6
3. Invisible opponent information	Scalar	Invisible information used for the value net	560
Opponent heroes	Position	Current 2D coordinates, distances, etc.	120
NPC	Position	Current 2D coordinates of all non-player characters, including minions, monsters, and turrets	440
4. Meta-Command feature	Spatial	2D image-like, extracted in channels for convolution	5x12x12
Meta-Commands	Spatial	All received Meta-Commands in the team	5x12x12

76 we adjusted the capability of agents to match the performance of strong humans by increasing the
77 observation delay (from 133ms to 200ms) and response delay (from 60ms to 120 ms).

78 A.9.2 Test Introduction

79 All participants were told the following instructions before testing:

- 80 • You will be invited into matches where your opponents and teammates are agents.
- 81 • Your goal is to win the game as much as possible by collaborating with agent teammates.
- 82 • You can collaborate with agent teammates through the in-game signaling system, just like
83 playing with human teammates.
- 84 • In addition, agent teammates will also send you signals representing their macro-strategies,
85 and you can judge whether to execute them based on your value system.
- 86 • Each game is about 10-20 minutes. Your identity information will not be disclosed to
87 anyone, and all game statistics are only used for academic research. You will voluntarily
88 choose whether to take the test.

89 If the participant volunteers to take the test, we will provide the equipment and game account to him,
90 and the test will begin.

91 A.9.3 Potential Participant Risks

92 First, we analyze the risks of this experiment to the participants. The potential participant risks of the
93 experiment mainly include the leakage of identity information and the time cost. And we have taken
94 a series of measures to prevent these risks.

95 Regarding identity information risks, our measures are as follows:

- 96 • We make a risk statement for participants and sign an identity information confidentiality
97 agreement.

- 98 • We only use game statistics without identity information in our research .
 - 99 • Special equipment and game accounts are provided to the participants to prevent leakage of
 - 100 equipment and account information.
 - 101 • The identity information of all participants is not disclosed to the public.
- 102 To compensate participants for their time cost, we offered each participant \$5 per match. Each match
- 103 is about 10-20 minutes, and participants can get about an average of \$20 an hour.
- 104 Finally, we have performed a process similar to IRB before the test is conducted. Our institution and
- 105 all participants have approved our research.

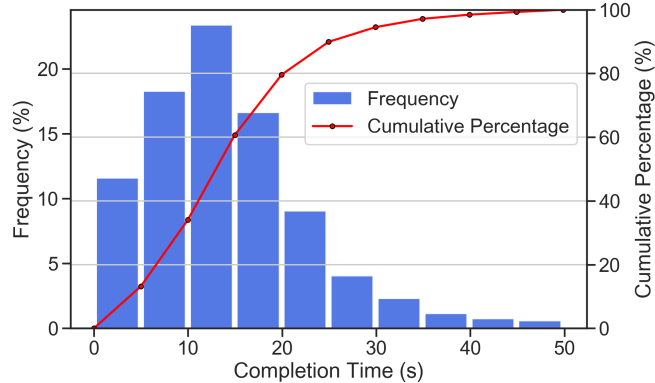


Figure 6: Time statistics of humans completing meta-commands in real games.

106 A.10 Additional Experimental Results

107 A.10.1 CEN

108 **Training Data.** We extract meta-commands from expert game replay authorized by the game provider,

109 which consist of high-level (top 1% player) license game data without identity information. The

110 input features of CEN are shown in Table 4. The game replay consists of multiple frames, and the

111 information of each frame is shown in Figure 1. For setting T^{mc} , we counted the player’s completion

112 time for meta-commands from expert game replay, and the results are shown in Fig. 6. We can

113 see that 80% meta-commands can be completed within the time of 20 seconds in *Honor of Kings*.

114 Thus, T^{mc} is set to 300 time steps (20 seconds). Given a state s_t in the trajectory, we first extract

115 the observation o_t for each hero. Then, we use a hand-crafted command extraction function f^{ce} to

116 extract the meta-command $m_t = f^{ce}(s_{t+T^{mc}})$ corresponding to the current state s_t in the future.

117 By setting up labels in this way, we expect the CEN $\pi_\phi(m|o)$ to learn the mapping from the current

118 observation o_t to its meta-command m_t . The detailed training data extraction process is as follows:

- 119 • First, we extract the trajectory $\tau = (s_0, \dots, s_t, \dots, s_{t+T^{mc}}, \dots, s_N)$ from the game replay, where
- 120 N is the total number of frames.
- 121 • Second, we randomly sample some frames $\{t|t \in \{0, 1, \dots, N\}\}$ from the trajectory τ .
- 122 • Third, for each frame t , we extract feature o_t from state s_t .
- 123 • Fourth, we extract the label m_t from the state $s_{t+T^{mc}}$ in frame $t + T^{mc}$, i.e. describe the state
- 124 using the meta-command space M .
- 125 • Finally, $\langle o_t, m_t \rangle$ is formed into a training pair as a sample in the training data.

126 **Optimization Objective.** After obtaining the dataset $\{\langle o, m \rangle\}$, we train the CEN $\pi_\phi(m|o)$

127 via supervised learning (SL). Due to the imbalance of samples at different locations of the meta-

128 commands, we use the focal loss [3] to alleviate this problem. Thus, the optimization objective

129 is:

$$L^{SL}(\phi) = \mathbb{E}_{O, M} \left[-\alpha m (1 - \pi_\phi(o))^\gamma \log(\pi_\phi(o)) - (1 - \alpha) (1 - m) \pi_\phi(o)^\gamma \log(1 - \pi_\phi(o)) \right],$$

130 where $\alpha = 0.75$ is the balanced weighting factor for positive class ($m = 1$) and $\gamma = 2$ is the tunable
 131 focusing parameter. Adam[2] is adopted as the optimizer with an initial learning rate of 0.0001.

132 **Experimental Results.** Figure 7 shows the meta-command distributions of the initial CEN, the
 133 converged CEN, and strong humans. We see that the meta-commands predicted by the CEN gradually
 134 converge from chaos to the meta-commands with important positions. And the distribution of the
 135 converged CEN in Figure 7(b) is close to the distribution of strong humans in Figure 7(c) and the
 136 corresponding KL divergence is 0.44, suggesting that the CEN can simulate the generation of human
 137 meta-commands in real games.

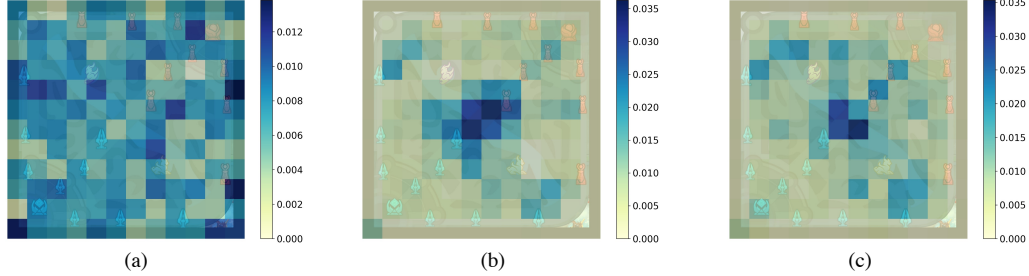


Figure 7: **The meta-command distributions of CEN and strong humans.** (a) The meta-command distribution of the initial CEN. (b) The meta-command distribution of the converged CEN. (c) The meta-command distribution of strong humans

138 A.10.2 MCCAN

139 **Optimization Objective.** The MCCAN is trained via goal-conditioned RL with the goal of achieving
 140 a near-human completion rate for the meta-commands generated by the pre-trained CEN while
 141 ensuring that the win rate is not reduced. We adopt an intrinsic reward to guide the process of
 142 executing the meta-command m_t :

$$r_t^{int}(s_t, m_t, s_{t+1}) = |f^{ce}(s_t) - m_t| - |f^{ce}(s_{t+1}) - m_t|,$$

143 where f^{ce} is a hand-crafted command extraction function. We train the MCCAN with the objective
 144 of maximizing the expectation over extrinsic and intrinsic discounted total rewards:

$$G_t = \mathbb{E}_{s \sim d_{\pi_\theta}, a \sim \pi_\theta} \left[\sum_{i=0}^{\infty} \gamma^i r_{t+i} + \alpha \sum_{j=0}^{T^{mc}} \gamma^j r_{t+j}^{int} \right],$$

145 where α is a trade-off parameter and $d_\pi(s) = \lim_{t \rightarrow \infty} P(s_t = s | s_0, \pi)$ is the probability when
 146 following π for t steps from s_0 .

147 **Training Process.** The MCCAN is trained by finetuning a pre-trained micro-action network [5]
 148 conditioned on the meta-command sampled from the pre-trained CEN. We modified the Dual-clip
 149 PPO algorithm [5] to introduce the meta-command m into the policy $\pi_\theta(a_t | o_t, m_t)$ and the advantage
 150 estimation $A_t = A(a_t, o_t, m_t)$. The Dual-clip PPO algorithm introduces another clipping parameter
 151 c to construct a lower bound for $r_t(\theta) = \frac{\pi_\theta(a_t | o_t, m_t)}{\pi_{\theta_{old}}(a_t | o_t, m_t)}$ when $A_t < 0$ and $r_t(\theta) \gg 0$. Thus, the
 152 policy loss is:

$$L^\pi(\theta) = \mathbb{E}_{s, m, a} [\max(cA_t, \min(\text{clip}(r_t(\theta), 1 - \tau, 1 + \tau)A_t, r_t(\theta)A_t)],$$

153 where τ is the original clip parameter in PPO. And the multi-head value loss is:

$$L^V(\theta) = \mathbb{E}_{s, m} \left[\sum_{head_k} (G_t^k - V_\theta^k(o_t, m_t)) \right], V_{total} = \sum_{head_k} w_k V_\theta^k(o_t, m_t),$$

154 where w_k is the weight of the k -th head and $V_\theta^k(o_t, m_t)$ is the k -th value.

155 **Experimental Results.** We conducted experiments to explore the influence of the extrinsic and
 156 intrinsic reward trade-off parameter α on the performance of MCCAN, and the win rate and comple-
 157 tion rate results are shown in Figure 8. We see that as α increase, the completion rate of MCCAN

158 gradually increases, and the winning rate of MCCAN first increases and then decreases rapidly. When
 159 $\alpha = 16$, the completion rate of the trained agent for meta-commands is 82%, which is close to the
 160 completion rate of humans (80%). And the win rate of the trained agent against the SOTA agent [1, 5]
 161 is close to 50%. Thus, we finally set $\alpha = 16$ in subsequent experiments.

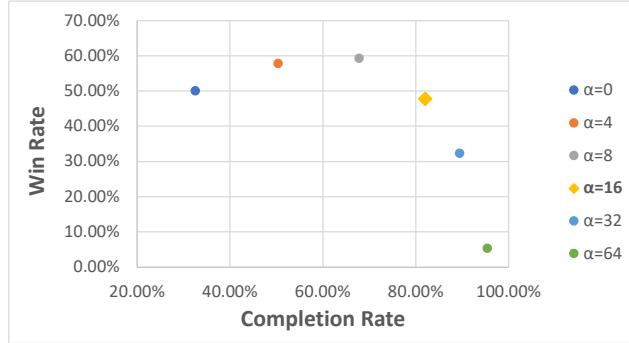


Figure 8: The win rate and completion rate of MCCAN with different α . The opponent is the pre-trained action network, i.e., MCCAN with $\alpha=0$.

Table 7: The WRs of different strong human-AI teams against the MC-Base agents in m AI + n Human mode.

Team Mode	Type of Agent		
	MC-Base	MC-Rand	MCC
2 AI + 3 Human	8%	3%	18%
3 AI + 2 Human	26%	18%	39%
4 AI + 1 Human	42%	28%	54%

162 A.10.3 Human-AI Game Test

163 **m AI + n Human Mode Result.** In addition to validating the generalization of the MCC agents
 164 to different levels of human teammates, we also evaluated the performance of the MCC agents to
 165 different numbers of human teammates. We had different numbers of strong humans team up with
 166 different types of agent teammates in m AI + n Human mode, where $m + n = 5$. We tested three
 167 team modes, including 2 AI + 3 Human mode, 3 AI + 2 Human mode, and 4 AI + 1 Human mode.
 168 The corresponding WR results are shown in Table 7. We can see that as the number of humans
 169 increases, the WR of the MC-Base-Human team drops dramatically as expected. Note that the WR
 170 of the SOTA [1, 5] agent-only team against the human-only team is close to 100% and the WR
 171 of the MC-Base agent against the SOTA agent is close to 50%(see in Fig. 8) Fortunately, when
 172 humans team up with MCC agents, they can achieve effective communication and collaboration
 173 on macro-strategies, resulting in significant increased WRs. We can also see that when humans
 174 team up with MC-Rand agents, the WR is the lowest, suggesting that randomly communicating and
 175 collaborating can greatly hurt performance.

Table 8: The subjective preference results of all participants in the Human-AI Game Test.

Participant Preference Metrics (from poor to perfect, 1~5)	Teammate	Type of Agent		
		MC-Base	MC-Rand	MCC
Reasonableness of H2A	Average Human	2.3 ± 0.38	2.7 ± 0.24	4.0 ± 0.6
	Strong Human	2.2 ± 0.21	2.5 ± 0.41	4.1 ± 0.55
Reasonableness of A2H	Average Human	-	1.9 ± 0.35	4.3 ± 0.31
	Strong Human	-	1.7 ± 0.24	4.4 ± 0.35
Overall Preference	Average Human	2.7 ± 0.41	1.3 ± 0.27	4.3 ± 0.4
	Strong Human	2.5 ± 0.21	1.2 ± 0.17	4.5 ± 0.41

176 **Subjective Preference Results.** During the Human-AI Game Test, after completing each game test,
177 the testers gave scores on several subjective preference metrics to evaluate their agent teammates,
178 including the Reasonableness of H2A (how well agents respond to the meta-commands sent from
179 testers), the Reasonableness of A2H (how reasonable the meta-commands sent from agents), and
180 the Overall Preference for agent teammates. We separate the scores of strong humans and average
181 humans and present the results in Table 8. We can see that for the Reasonableness of H2A metric,
182 both strong and average humans gave the highest scores to MCC agents, which are significantly
183 higher than that of other agents, indicating that humans relatively agree with the value estimation of
184 MCC agents on meta-commands sent from humans. This is also verified in Fig. 7 of the main text.
185 We can also see that for the Reasonableness of A2H metric, humans rated MCC agents much better
186 than MC-Rand agents, indicating that humans believe that the meta-commands sent from MCC agents
187 are more aligned with their own value system, so humans are more willing to trust and collaborate
188 with MCC agents. For the Overall Preference metric, humans are satisfied with teaming up with
189 MCC agents, scoring the highest scores compared to other agents. The results of these subjective
190 preference metrics are also consistent with the results of objective metrics (Tables 1,2 of main text
191 and Table 7 of Appendix).

192 **A.11 Limitations and Future work**

193 **A.11.1 Limitations**

194 There are three main limitations to our research work. 1) Due to the complexity of the MOBA game
195 and the complexity of the MCC framework, the MCC framework adopts a sequential training manner
196 instead of an end-to-end training manner. Thus, the training process of the MCC framework is tedious.
197 2) The training of the MCC agent consumes a lot of computing resources like the training of the
198 SOTA MOBA AI agent. Thus, the computational cost of extending the MCC framework to other
199 complex MOBA games is huge. 3) The meta-command we proposed is generic only to MOBA games
200 and cannot be directly extended to other types of games, such as First-Person Shooting (FPS) and
201 Massively Multiplayer Online (MMO).

202 **A.11.2 Future work**

203 **From the application side,** we will precipitate this research work and apply it to the friendly bots in
204 teaching mode of *Honor of Kings*, aiming to provide gameplay teaching to novice players.

205 **From the research side,** first of all, we will optimize the training process of the MCC framework,
206 including the training process of the SOTA AI systems, reduce the computing resources required for
207 training the MOBA agent, aiming to lower the threshold for researchers to study and reproduce work
208 on MOBA games. Second, we will design a more general meta-command representation, such as
209 natural language, and extend the MCC framework to other types of games. All in all, it is our sincere
210 hope that human-AI collaboration in complex environments will attract more and more researchers'
211 attention, and we also hope that this work can provide researchers with some new ideas.

212 **References**

- 213 [1] Yiming Gao, Bei Shi, Xueying Du, Liang Wang, Guangwei Chen, Zhenjie Lian, Fuhao Qiu,
214 Guoan Han, Weixuan Wang, Deheng Ye, et al. Learning diverse policies in moba games via
215 macro-goals. *Advances in Neural Information Processing Systems*, 34, 2021.
- 216 [2] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
217 *arXiv:1412.6980*, 2014.
- 218 [3] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense
219 object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages
220 2980–2988, 2017.
- 221 [4] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster,
222 and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent
223 reinforcement learning. In *Proceedings of International Conference on Machine Learning*, pages
224 4295–4304, 2018.

- 225 [5] Deheng Ye, Guibin Chen, Wen Zhang, Sheng Chen, Bo Yuan, Bo Liu, Jia Chen, Zhao Liu, Fuhao
226 Qiu, Hongsheng Yu, et al. Towards playing full moba games with deep reinforcement learning.
227 *Advances in Neural Information Processing Systems*, 33:621–632, 2020.