

376 A Author Statement

377 We have released HoK3v3 - the Honor of Kings 3v3 Arena - as an open-source project under
378 the Apache License V2.0. The relevant code and gamecore can be found at https://github.com/tencent-ailab/hok_env, while detailed documentation for HoK3v3 is available at https://doc.aiarena.tencent.com/paper/hok3v3/latest/hok3v3_env/honor-of-kings/. All
379 experiments can be replicated using the source code, inclusive of hyperparameters and configurations.
380 The game developers of Honor of Kings have authorized HoK3v3, and the authors bear all responsibility
381 in case of violation of rights. We will also ensure data accessibility and provide necessary
382 maintenance.
383
384

385 B Environment Details

386 B.1 Mechanisms

387 There are two types of mechanisms present in HoK3v3: *Crystal* and *Turret*. Moreover, the Arena
388 contains two *Turrets* and one *Crystal*. The primary objective for players is to engage in combat to
389 safeguard their own *Crystal* while simultaneously attempting to destroy the opposing team’s *Crystal*.
390 The *Crystal* possesses a 90% resistance to damage, and it appears in the map within 10 seconds of
391 the game’s initiation.

392 *Turrets* are classified into two types: the *Vanguard Turret* and the *Base Turret*. Once the *Vanguard*
393 *Turret* is destroyed, players can proceed to destroy the *Base Turret*, followed by the final *Crystal*.
394 Both the *Vanguard Turret* and the *Base Turret* are formed 10 seconds after the game begins. For
395 the initial 2 minutes, a protection mechanism is in place to safeguard the *Turrets*, allowing them
396 to withstand 80 points of normal attack damage from heroes. Additionally, *Turrets* enjoy a 55%
397 damage-free rate in the absence of *Creeps*.

398 Sixty seconds after the commencement of the game, an HP pack will become available behind both
399 the *Vanguard Turret* and the *Base Turret*, providing heroes with the means to restore their health
400 points (HP). Once utilized, the HP pack will disappear and subsequently replenish itself every 60
401 seconds. However, in the event that either *Turret* is destroyed, the regeneration of the HP pack of this
402 *Turret* itself will cease.

403 The defense attributes of the mechanisms are shown in Table. 3.

Table 3: Defense attributes of the mechanisms.

Mechanisms	Basic HP	Growth HP	Basic Armor	Basic Resistance
<i>Vanguard Turret</i>	6000	700	200	200
<i>Base Turret</i>	10000	700	200	200
<i>Crystal</i>	8000	600	200	200

404 The mechanisms select attack targets according to a consistent rule. If an enemy hero fails to inflict
405 damage on ally heroes within the *Turret*’s attack range, the *Turret* will prioritize attacking the first
406 unit that enters its range. Once the initial unit is eliminated, the *Turret* will then proceed to attack
407 minions, summoned creatures, and heroes, in that order. In cases where units share the same priority,
408 the *Turret* will direct its attacks towards the nearest unit. However, when an enemy hero inflicts
409 damage on ally heroes, the *Turret* will focus on the first enemy hero responsible for the damage. This
410 targeting persists until the enemy hero either exits the attack range or is eliminated.

411 The extent of damage inflicted by the mechanisms will accumulate with each subsequent attack, and
412 this damage is characterized as physical damage that bypasses any defensive measures. The attacks
413 performed by these mechanisms are listed in Table. 4.

414 Destroying enemy mechanisms will gain experience and golds, as shown in Table. 4.

Table 4: Attacks performed by mechanisms.

Mechanisms	Basic Attack	Attack Bonus	Maximum Attack Bonus	Experience	Gold
<i>Vanguard Turret</i>	430	300	1500	100	100
<i>Base Turret</i>	500	300	1500	100	100
<i>Crystal</i>	580	300	3000	0	0

415 B.2 Creep

416 The creep serves as the primary source of experience and gold in the HoK3v3, constituting the largest
 417 proportion of these resources within the entire game. It is primarily categorized into two types:
 418 ordinary creep and super creep. Upon the player's destruction of the *Base Turret*, the ordinary
 419 creep is substituted with the more formidable super creep. Creep materializes 12 seconds after the
 420 commencement of the game and subsequently regenerate every 24 seconds.

421 In the first four minutes of the game, the composition of the creep consists of two *Warriors* and two
 422 *Mages* respectively. As the game progresses, the creep composition changes to include two *Warriors*,
 423 one *Mage*, and one *Catapult*. Upon the player's destruction of the *Base Turret*, the *Ordinary Creep* is
 424 replaced by the more formidable *Super Creep*, which is comprised of two *Warriors*, one *Mage*, and
 425 one *Super Warrior*. The essential attributes of the creep are presented in Table 5.

Table 5: Basic attributes of creep.

Creep	Attack	Magic	Armor	Resistance	HP	Experience	Gold
<i>Warrior</i>	60	60	0	0	1860	60	48
<i>Mage</i>	120	120	0	0	1545	45	36
<i>Catapult</i>	192	192	0	0	2790	100	84
<i>Super Warrior</i>	576	576	183	183	4185	100	70

426 B.3 Jungles

427 The jungles serve as the primary source of experience and gold for *Assassin Heroes* such as "*Zhaoyun*".
 428 In this particular map, both sides' players have access to the entire jungle area. The jungles contain
 429 various types of monsters, including normal creatures, the formidable "*Scarlet Statue*," the elusive
 430 "*Treasure Thief*", and the powerful "*Tyrant*". The spatial distribution of various monster species
 431 within the jungles can be referenced in Figure. 7. Upon slaying the "*Dark Wolf*," a hero receives the
 432 "*Forest's Roar*" buff, granting them a 30% increase in movement speed and a 20% reduction in skill
 433 cooldowns for a duration of 30 seconds. However, this buff dissipates upon the hero's death.

434 Defeating the "*Tyrant*" bestows the hero with the "*Tyrant's Power*." Accumulating multiple "*Tyrant*"
 435 kills further enhances this power. The effects of slaying the "*Tyrant*" multiple times are as follows:
 436 the first kill increases the HP recovery of all allies by 1% every 2 seconds, the second kill augments
 437 the damage inflicted by all allies against enemy mechanisms, and the third kill amplifies the physical
 438 and magical output of all allies by 30%. For additional information regarding the specifics of the
 439 jungles, please refer to Table 6.

Table 6: Details about the jungles.

Creature Name	Attack	Magic	Armor	Resistance	HP	Experience	Gold
<i>Big / Little Demon Vanguard</i>	138 82	138 82	183 109	183 109	2480 1488	60 30	60 30
<i>Big / Little Archer</i>	138 82	138 82	183 109	183 109	2480 1488	60 30	60 30
<i>Big / Little White-tail Deer</i>	138 82	138 82	183 109	183 109	2480 1488	60 30	60 30
<i>Big / Little Dark Wolf</i>	138 82	138 82	183 109	183 109	2480 1488	70 50	70 50
<i>Scarlet Statue</i>	216	216	183	183	3720	90	90
<i>Treasure Thief</i>	204	204	183	183	5400	160	135
<i>Tyrant</i>	204	204	183	183	9000	300	200

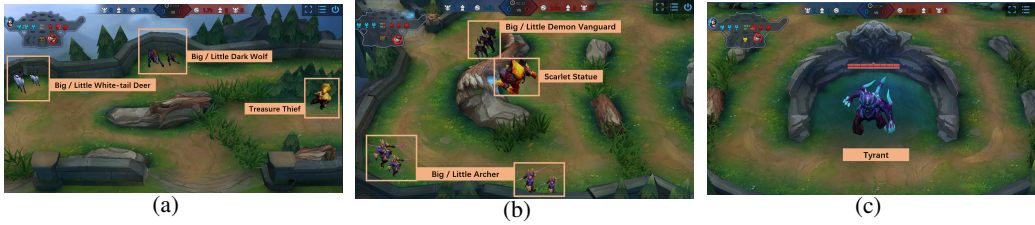


Figure 7: The spatial distribution of various monster species within the jungles.

440 B.4 Heroes

441 In the HoK3v3, we have open-sourced a total of 30 heroes (Fig. 8), which can be classified into three
 442 types: *Mage*, *Marksmen*, and *Assassin*. Each type consists of 10 heroes. *Mage* and *Marksmen* heroes
 443 primarily operate in the middle lane, acquiring experience and gold by eliminating opponent heroes
 444 or creeps. On the other hand, *Assassin* heroes predominantly operate in the jungle, killing monsters
 445 to obtain gold and experience. Additionally, *Assassin* heroes also venture into the middle lane to col-
 446 laborate with *Mage* and *Marksmen* heroes in eliminating opponent heroes or destroying *Turrets*. Each
 447 agent is able to select one hero to control and cooperate with the other two agents. For further details
 448 regarding hero skills, please consult the website: <https://pvp.qq.com/m/m201706/heroList.shtml>.

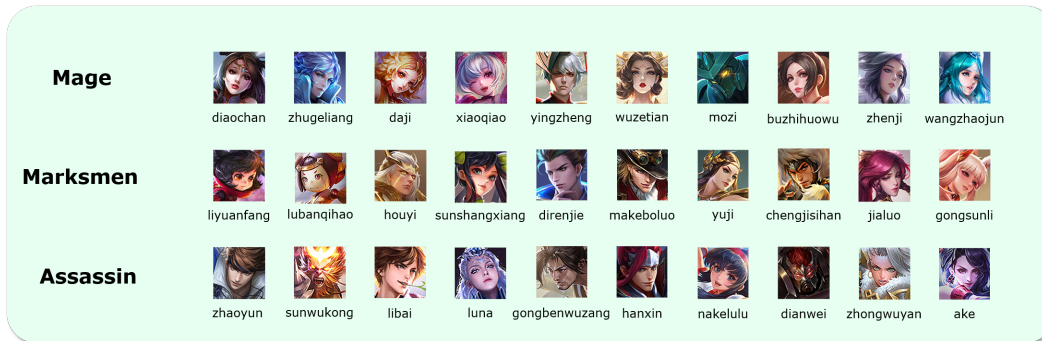


Figure 8: Details of 30 open-sourced heroes.

449 C Agent Details

450 C.1 Observation Space

451 In order to enhance the training process and facilitate academic exploration, the HoK3v3 incorporates
 452 various components, including the intricate game interaction logic, the training framework project,
 453 and the feature design project. By encapsulating these elements, the HoK3v3 offers comprehensive
 454 and essential information. The complete observation space dimension comprises 4586, and Table 7
 455 provides a detailed breakdown of feature categories, descriptions, and dimensions.

456 Note:

- 457 • All the features, except *FeatureWholeInfo*, consist of both absolute information features and
 458 relative information features. Let us consider *FeatureHero* as an illustrative example. The absolute
 459 information features, such as hero ID, blood, HP, and attack power, remain consistent for all three
 460 heroes. On the other hand, the relative information features, such as the x-axis, z-axis, distance, and
 461 other dimensions pertaining to the current player, vary among the three heroes.
- 462 • The units related to the two sides are presented in the following order: our side and enemy side.
 463 For example, the *FeatureHero* unit represents [3 heroes on our side, 3 heroes on the enemy side],

464 the *FeatureSoldier* unit represents [10 soldiers on our side, 10 soldiers on the enemy side], and the
 465 *FeatureOrgan* unit represents [3 turrets on our side, 3 turrets on the enemy side].

466 • Units that are not associated with specific sides are arranged according to the ID of the unit. For
 467 example, the *FeatureMonster* unit consists of [*Monster 0*, *Monster 1*, ..., *Monster N*], with the *Tyrant*
 468 unit being positioned last within the *Monster* unit.

Table 7: The categories of features, descriptions, and the dimensions.

Categories of Features	Descriptions	Dimensions	Start Index	End Index
<i>FeatureImg</i>	Image-like features, including 6 channels such as obstacle channel and grass channel.	6*17*17	0	1733
<i>FeatureHero</i>	From vision of the current player, state information of 6 heroes from both sides, i.e. hero ID, HP, etc.	6*251	1734	3239
<i>FeatureMainHero</i>	Private features of the current hero, i.e. whether the enemy hero is within the attack range of the current hero.	44	3240	3283
<i>FeatureSoldier</i>	The state of 20 <i>Creep</i> of allies and enemies: types, HP, positions, etc.	20*25	3284	3783
<i>FeatureOrgan</i>	The state of 6 <i>Turrets</i> : types, HP, positions, etc.	6*29	3784	3957
<i>FeatureMonster</i>	The state of 20 <i>Monsters</i> : types, HP, positions, etc.	20*28	3958	4517
<i>FeatureWholeInfo</i>	Gold of allies and enemies; kills, surviving turrets, etc.	68	4518	4585

469 C.2 Action Space

470 The original action space in the HoK3v3 comprises a triad of action buttons: the move direction button,
 471 the skill offsets on the x- and z-axes button, and the target game units button. This comprehensive set
 472 encompasses all possible actions that the hero can undertake in a hierarchical fashion.

473 Specifically, the player must make decisions regarding the following aspects:

474 **Selection of action button:** The player needs to determine which action button to choose, such as
 475 the move button, attack button, skill button, return button, and so on.

476 **Execution details:** The player must specify the precise execution details, including controlling the
 477 direction of movement and managing the landing position of skills.

478 **Target selection:** The player must decide which target to select for the intended action.

479 Details of the action space can be referred to Table. 8.

480 C.3 Legal Action Mask

481 As shown in Table 8, in each time step of an episode, every hero has the option to choose one action.
 482 However, their choice of action is not arbitrary. Therefore, a legal action mask exists for each time
 483 step, restricting the hero from selecting illegal actions. The dimension of the legal action mask is the
 484 same as that of the action when it comes to the action types "*Button*", "*Move*", and "*Skill*". However,
 485 there is a dependency relationship in the legal action mask between the action types "*Target*" and
 486 "*Button*": the legal actions of "*Target*" depend on the chosen "*Button*". In other words, only when a
 487 hero chooses a "*Button*" can they determine the legal actions for "*Target*". Since there are a total of
 488 13 "*Buttons*" and 7 "*Targets*", the dimension of the legal action mask for "*Target*" is 13×7 .

Table 8: Details of the action space.

Action Type	Sub Action	Description	Dimension
Button	None	inactive state	1
	None	inactive state	1
	Move	move hero	1
	Normal Attack	cast normal attack	1
	Skill 1	cast skill 1	1
	Skill 2	cast skill 2	1
	Skill 3	cast skill 3	1
	Skill 4	cast skill 4 (for specific heroes)	1
	Chosen Skill	cast chosen skill	1
	Recall	return to the base	1
	Equipment Skill	cast equipment skill	1
	Heal Skill	cast heal skill	1
	Friend Skill	cast friend skill (for specific heroes)	1
Move	Move Dir	move direction	25
Skill	Skill X	skill offsets on the x-axis	42
	Skill Z	skill offsets on the z-axis	42
Target	None	no target	1
	Enemy	3 enemy heroes	3
	Friend	3 friend heroes	3
	Self	own hero	1
	Monster	20 monsters	20
	Soldier	10 closest monsters	10
	Turret	the closest turret	1

489 C.4 Sub Action Mask

490 Action masking refers to the process of removing certain actions that cannot be executed simul-
 491 taneously with the current action, resulting in a selection of permissible actions. To facilitate
 492 comprehension, we present a couple of examples:

493 **Example-1:** Upon choosing the Button-Move action, only the Move Dir sub-actions remain after the
 494 mask is applied. These sub-actions enable control over the direction of movement, as depicted in
 495 Fig. 9.

496 **Example-2:** Upon selecting the Button-Normal Attack action, only the Target sub-action remains,
 497 defining the target of the normal attack after the mask is applied.

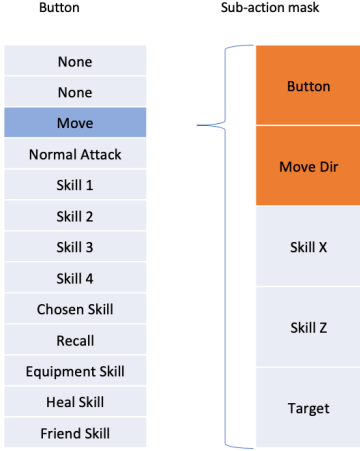
498 Similar masking principles apply to other actions as well. It is important to note that the specific
 499 sub-action masks may vary depending on the heroes and equipment involved. For further information,
 500 please consult the official website at "<https://pvp.qq.com/web201605/herolist.shtml>".

501 C.5 Reward

502 The design of rewards takes the following factors into consideration:

- 503 • Hero Development: The rewards are based on the golds and experience gained through killing
 504 monsters and creeps.
- 505 • KDA (Kills, Deaths, and Assists): The rewards are influenced by the player's performance in
 506 terms of kills, deaths, and assists.
- 507 • Hero's State: The rewards are tied to the hero's remaining HP (Health Points).
- 508 • Game Progression: The rewards are determined by the destruction of Turrets and the Crystal.

For specific details regarding the design of rewards, please refer to Table 10.



Reward Factors	Descriptions	Weights
<i>hp_rate_sqrt_sqrt</i>	the fourth root of Hp rate	1.0
<i>money</i>	economic growth	0.001
<i>exp</i>	experience growth	0.001
<i>tower</i>	HP of own <i>Turret</i>	1.0
<i>killCnt</i>	kill count	1.0
<i>DeadCnt</i>	dead count	-1.0
<i>assistCnt</i>	assist count	1.0
<i>total_hurt_to_hero</i>	damage to heroes	0.1
<i>atk_monster</i>	attack monster	0.1
<i>win_crystal</i>	destroy enemy crystal	1.0
<i>atk_crystal</i>	attack enemy crystal	1.0

Figure 10: The specific design of rewards.

Figure 9: Sub Action mask after selecting the Button-Move action.

509

510 When it comes to calculating the final reward, we employ a "zero-sum" approach wherein the average
 511 reward of the enemy camp is subtracted. The individual hero reward is determined by applying
 512 weights, which can be referenced in Table 10. As demonstrated in Equation 1, the final reward is
 513 obtained by utilizing each hero’s $hero_{reward_{zero_sum}}$ value.

$$\begin{aligned}
 hero_{reward} &= w_1 r_1 + w_2 r_2 + \dots + w_n r_n \\
 team_{reward} &= \frac{1}{3} \sum_{i=1}^3 hero_{reward} \\
 hero_{reward_{zero_sum}} &= \begin{cases} hero_{reward} - team_{reward_{camp2}}, & \text{if hero in camp1} \\ hero_{reward} - team_{reward_{camp1}}, & \text{if hero in camp2} \end{cases}
 \end{aligned} \tag{1}$$

514 **C.6 Comparison with related works**

Table 9: A detailed comparison of HoK3v3 and other related works.

	Observation Space	Action Space	Focus	Heterogeneity [†]
Google Research Football	115*	19	Cooperation	✗
StarCraft Multi-Agent Challenge	16-336	7-70	Cooperation	o
Multi-agent MuJoCo	≤376	≤17	Cooperation	o
Multi-agent Particle Environment	≤20**	≤10**	Cooperation&Competition	✗
Melting Pot	88x88x3	6***	Cooperation&Competition	✗
HoK Arena	491	83	Competition	✓
HoK3v3 (Ours)	4586	161	Cooperation&Competition	✓

[†] ✗ symbolizes homogeneity, "o" indicates only numerical heterogeneity, and ✓ denotes true heterogeneity.

* The floats representation proposed in the original paper.

** Estimations for common scenarios.

*** Common movement actions proposed in original paper.

515 **D Elo Details**

516 The Elo rating system is a common ranking system used in competitive matches[2]. Here is the
 517 calculation method of the Elo rating system used in this article:

518 1. Each model has an Elo rating, which represents their skill level in the competition. The initial Elo
519 rating of each model is 1500.

520 2. In each match, the expected probability of winning for each model is calculated based on their Elo
521 rating difference. The formula for expected probability is:

$$E_a = \frac{1}{1 + 10^{\left(\frac{R_b - R_a}{400}\right)}} \quad (2)$$

522

$$E_b = 1 - E_a \quad (3)$$

523 where ' R_a ' and ' R_b ' are the ELO ratings of the two models.

524 3. After each match, the Elo ratings of the two models are updated based on the actual result. If
525 model A wins, its actual score S_a is 1 and model B's actual score is 0. The formula for updating Elo
526 rating is:

527

$$R'_a = E_a + K \times (S_a - E_a) \quad (4)$$

$$R'_b = R_b + K \times ((1 - S_a) - E_b) \quad (5)$$

528 where K is a constant that determines the amount of change in Elo rating after each match. In this
529 article, $K = 40$.

530 4. In this article, each model plays 128 matches against other models, and the win-loss records are
531 shuffled randomly. The Elo ratings are updated based on the shuffled win-loss records, and this step
532 is repeated for 200 times to calculate the average Elo ratings. This is done to reduce the error caused
533 by different match orders.

534 E Hyperparameters

535 We have included Table. 10, 11, 12 and 13, which present the key hyperparameters utilized in the
536 relevant experiments. This table encompasses the essential parameters required for conducting the
537 experiments effectively.

538 F Sub-tasks Details

539 In these sub-tasks, we modify the reward function to individual reward item that corresponds
540 specifically to the given sub-task, with weight as 1, in stead of weighted sum of multiple items. The
541 concrete details of sub-tasks are as follows:

542 • **Gold:** Obtaining more gold generally provides a significant advantage, as it is the most important
543 resource in the game. The objective of this sub-task is to collect more gold by destroying enemy units
544 (heroes, creeps, and turrets) or monsters. The corresponding reward item for this sub-task is 'money'.

545 • **Exp:** Similar to gold, experience points are crucial in the game as they determine the level of
546 the heroes. Therefore, we have designed this sub-task to modify the objective to focus on gaining
547 more experience points, which helps heroes level up faster. The corresponding reward item for this
548 sub-task is 'exp'.

549 • **Kill:** As a competitive game, killing an enemy provides both gold and experience points, while
550 also temporarily incapacitating the slain enemy, thus granting a significant advantage to the team.
551 Consequently, we have designed the **Kill** sub-task to specifically train the agents to eliminate enemies
552 as frequently as possible. The corresponding reward item for this sub-task is 'killCnt'.

553 • **Hurt:** The rate of hurting is another metric that signifies killing enemies, and it carries a more
554 dense reward. This sub-task serves as an alternative to the **Kill** objective, aiming to maximize the
555 extent of hurt inflicted. The corresponding reward item for this sub-task is 'total_hurt_to_hero'.

556 • **Turret:** Destroying the turrets of enemies is a crucial sub-goal in the game, as it grants access to
557 the crystal. We have designed the **Turret** sub-task to enhance players' abilities in destroying enemy
558 turrets and defending themselves. The corresponding reward item for this sub-task is 'tower'.

Table 10: Hyperparameters.

Hyperparameters	Value
Batch Size	288
γ	0.995
LSTM Time Steps	16
λ in GAE	0.95
PPO Clip ϵ	0.2
PPO Clip c	3.0
Optimizer	Adam
beta1	0.9
beta2	0.999
eps	1.00E-08
Learning Rate	6.00E-04

Table 11: Reward weight of "Mages".

Reward Factors	Weights
<i>hp_rate_sqrt_sqrt</i>	3.0
<i>money</i>	0.005
<i>exp</i>	0.0
<i>tower</i>	1.0
<i>killCnt</i>	1.0
<i>DeadCnt</i>	0.0
<i>assistCnt</i>	1.0
<i>total_hurt_to_hero</i>	0.3
<i>atk_monster</i>	0.0
<i>win_crystal</i>	0.0
<i>atk_crystal</i>	0.0

Table 12: Reward weight of "Marksmen".

Reward Factors	Weights
<i>hp_rate_sqrt_sqrt</i>	3.0
<i>money</i>	0.005
<i>exp</i>	0.0
<i>tower</i>	1.0
<i>killCnt</i>	1.0
<i>DeadCnt</i>	0.0
<i>assistCnt</i>	1.0
<i>total_hurt_to_hero</i>	0.3
<i>atk_monster</i>	0.0
<i>win_crystal</i>	0.0
<i>atk_crystal</i>	0.0

Table 13: Reward weight of "Assassins".

Reward Factors	Weights
<i>hp_rate_sqrt_sqrt</i>	3.0
<i>money</i>	0.005
<i>exp</i>	0.0
<i>tower</i>	1.0
<i>killCnt</i>	1.0
<i>DeadCnt</i>	0.0
<i>assistCnt</i>	1.0
<i>total_hurt_to_hero</i>	0.3
<i>atk_monster</i>	0.02
<i>win_crystal</i>	0.0
<i>atk_crystal</i>	0.0

559 • **Monster:** Monsters residing in the jungle play a vital role in enabling heroes to acquire gold
560 and experience points. Instead of the original objective, we have modified it to focus on attacking
561 monsters. The corresponding reward item for this sub-task is 'atk_monster'.

562 G Baseline Details

563 **Encoder:** We employ a meticulously designed backbone as our feature extractor, specifically tailored
564 to handle the extensive observation space. The encoder consists of multiple units dedicated to process-
565 ing different aspects of the observation. These units encompass a convolution module, responsible for
566 extracting image-like features, as well as distinct modules for hero, creep, turret, monster, and game
567 status information. Please refer to the code https://github.com/tencent-ailab/hok_env for
568 the implementation details of the network.

569 **PPO:** We employ Proximal Policy Optimization (PPO) [15] as our baseline algorithm. Specifically,
570 we adopt the dual-clip version of PPO, which has been empirically validated as effective in similar
571 environments [22, 20]. Additionally, to address the challenges posed by multi-agent control, we
572 incorporate the independent learning paradigm [17] and parameter sharing with PPO.

573 **CPPO:** CPPO is a communication-based variant of PPO. This approach shares similarities with the
574 standard implementation of PPO, with the exception that a portion of the processed feature from each
575 hero undergoes max-pooling to obtain a shared feature. This shared feature is then utilized by both
576 the policy and value networks.

577 **MAPPO:** To evaluate the efficacy of the CTDE paradigm, we employ a CTDE variant of PPO, known
578 as MAPPO [23]. In contrast to independent learning, MAPPO incorporates a unified global value

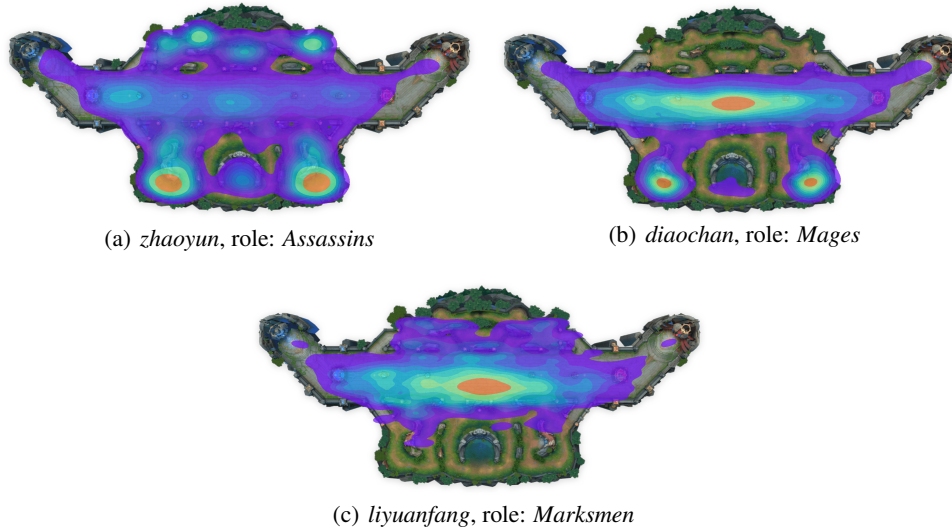


Figure 11: The trajectories of heroes with different roles.

579 network that can access the processed features of all heroes and utilizes the average reward of the
 580 heroes within a team as a shared reward.

581 H Additional Experiments

582 H.1 Visualizing Hero Trajectories

583 For the purpose of better illustrating the trajectories of heroes with different roles, namely *Assassins*,
 584 *Mages*, and *Marksmen*, we have constructed a heatmap utilizing the hero locations of 50 different
 585 trajectories. These trajectories were generated using the *Level-5 Model* with identical lineups on
 586 both sides, consisting of *zhaoyun*, *diaochan*, and *liyuanfang*, representing the *Assassins*, *Mages*, and
 587 *Marksmen* respectively. As depicted in Fig. 11, the heatmap is based on the HoK3v3 map.

588 When examining the *Assassins*, exemplified by *zhaoyun* (Fig. 11(a)), the heatmap illustrates the
 589 trajectory of *zhaoyun* as it permeates the entire map. This aligns with the role of *Assassins*, who
 590 are tasked with eliminating monsters throughout the jungle, acquiring gold and experience, and
 591 collaborating with *Mages* and *Marksmen* to eliminate opposing heroes. In the case of *Mages*,
 592 represented by *diaochan* (Fig. 11(b)), the heatmap reveals that *diaochan* predominantly remains within
 593 the middle lane, occasionally venturing into the lower section of the jungle to accrue resources by
 594 slaying monsters. As for *Marksmen*, embodied by *liyuanfang* (Fig. 11(c)), the heatmap demonstrates
 595 that *liyuanfang* primarily operates within the middle lane, but occasionally ventures into the upper
 596 region of the jungle to obtain gold and experience by dispatching monsters.

597 These findings indicate that *Mages*, such as *diaochan*, and *Marksmen*, like *liyuanfang*, strategically
 598 exploit resources in different areas of the jungle, specifically the lower and upper sections respectively,
 599 to mitigate competition for jungle resources. In summary, the map is effectively utilized by all three
 600 hero roles, each exhibiting distinctive characteristics.

601 H.2 Ablation Study on Whether to Use Shared Rewards

602 For the experiments conducted in our paper, we utilized the **PPO** and **CPPO** algorithms. In each
 603 lineup, consisting of three heroes from both sides, we employed separate rewards for training
 604 purposes. In other words, different heroes within each lineup received distinct rewards during training.
 605 Consequently, we conducted an ablation study using the **CPPO** algorithm to investigate the impact of
 606 utilizing averaged shared rewards during the training stage. The results are presented in Fig. 12(a),

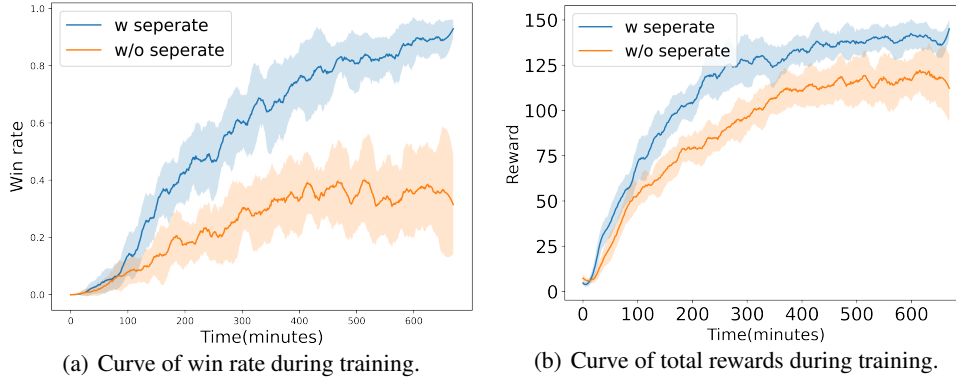


Figure 12: Ablation study on whether to use shared rewards during the training stage under three random seeds. **w seperate**: heroes in every lineup using different rewards. **w/o seperate**: heroes in every lineup using the same averaged rewards.

607 which depicts the winning rate curve during training with both shared and no_shared rewards. The
 608 figure clearly illustrates a significant performance drop when using shared rewards, which can be
 609 attributed to the diverse roles of heroes, such as "Marksmen" and "Mages", who possess distinct
 610 fighting styles. Therefore, it is appropriate to assign separate rewards to each hero based on their
 611 specific roles. Additionally, as shown in Fig. 12(b), the total rewards obtained by heroes when using
 612 shared rewards are substantially lower compared to those using no_shared rewards. Consequently,
 613 for all other experiments, we employ no_shared rewards.

614 H.3 Ablation Study on Whether to Use Zero-Sum Rewards

615 As described in Appendix. C.5, all the experiments conducted in our paper utilize zero-sum rewards
 616 to train, see Equation. 1. Therefore, we have conducted an ablation study to investigate the impact of
 617 using zero-sum rewards, as depicted in Fig.13. In other words, when employing zero-sum rewards,
 618 calculations are based on Equation. 1, whereas without zero-sum rewards, Equation. 6 is used.
 619 Fig. 13(a) demonstrates that during the later stages of training, the *w/o zero-sum* approach performs
 620 significantly worse than the **w zero-sum** approach, indicating that the utilization of zero-sum rewards
 621 can enhance the upper limit of the algorithm. Additionally, Fig. 13(b) illustrates that the total rewards
 622 obtained with zero-sum calculations are fewer compared to those obtained without zero-sum. This
 623 discrepancy arises from the subtraction of the average reward of the enemy camp (Equation. 1) when
 624 zero-sum rewards are employed. Consequently, for all the other experiments, we employ zero-sum
 625 rewards.

$$hero_{reward} = w_1r_1 + w_2r_2 + \dots + w_nr_n \quad (6)$$

626 H.4 Ablation Study on Whether to Use Sub-Action Mask

627 As described in Appendix C.4, sub-action mask is utilized to exclude certain *Action Types* that cannot
 628 be executed simultaneously with the current action. In other words, during each time step, not all
 629 *Action Types* (Table. 8), namely "Button", "Move", "Skill-X", "Skill-Z" and "Target" are necessary for
 630 training. Therefore, a sub-action mask is employed to eliminate the irrelevant *Action Types*. Therefore,
 631 an ablation study is conducted, as illustrated in Fig. 14, to examine the use of the sub-action mask.
 632 In Fig. 14(a), the winning rate of the *w/o SAM* (without sub-action mask) condition is considerably
 633 lower than that of the *w SAM* (with sub-action mask) condition. This discrepancy can be attributed to
 634 the fact that, during the backward stage of *w/o SAM*, gradients of the irrelevant *Action Types* can be
 635 regarded as noise interfering with the learning process of others. In Fig. 14(b), it can be observed

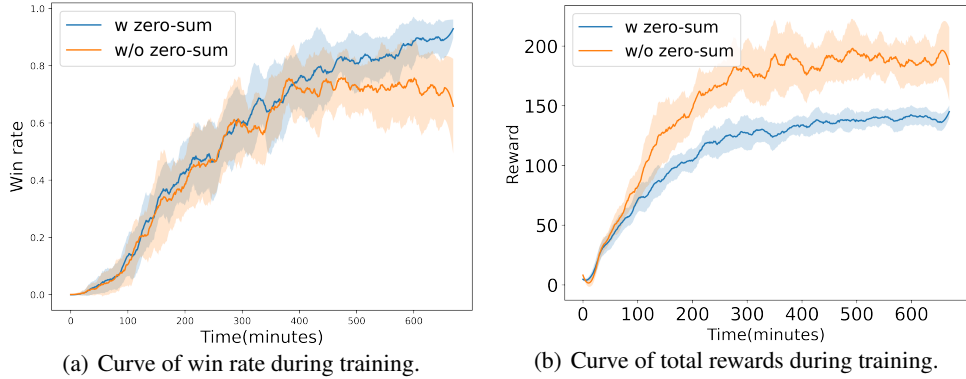


Figure 13: Ablation study on whether to use zero-sum rewards during the training stage under three random seeds. **w zero-sum**: Rewards are calculated by Equation. 1. **w/o zero-sum**: Rewards are calculated by Equation. 6.

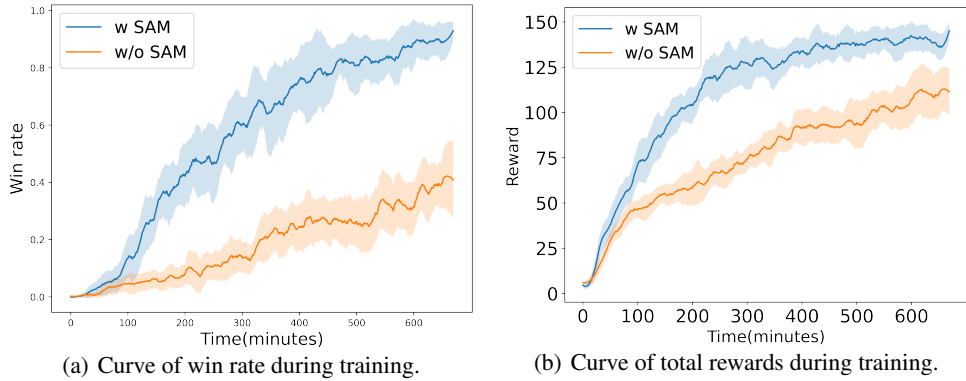


Figure 14: Ablation study on whether to use sub-action mask(SAM) during the training stage under three random seeds. **w SAM**: Using sub-action mask during training. **w/o SAM**: Without sub-action mask during training.

636 that the total rewards obtained in the *w/o SAM* condition are fewer than those achieved in the *w SAM*
 637 condition, as expected. Consequently, for all the other experiments, we employ the sub-action mask.

638 I Limitations and Future Works

639 The complexity and realism of Honor of Kings provide more opportunities for diverse research
 640 directions, which have not been explored thoroughly and are crucial for our future work. We encourage
 641 broader community involvement in studying this environment. Besides, We plan to optimize the
 642 deployment of our environment to fit for multiple platforms and organize more competitions based
 643 on the Honor of Kings to expand the influence of the Honor of Kings environment and encourage
 644 research on reinforcement learning.

645 J Additional Discussion

646 J.1 Discussion on the meaning os sub-tasks

647 Each individual subtask can be conceptualized as a sub-goal, representing a breakdown of the overall
 648 objective. For instance, the objective of destroying the enemies' crystal can be roughly broken down
 649 into sub-goals such as acquiring gold and experience points by defeating monsters -> hurt and kill

650 the enemies -> destroying their turrets -> ultimately shattering their crystal. This decomposition of
651 the main goal has the potential to significantly enhance hierarchical and goal-driven research in multi-
652 agent reinforcement learning (MARL). Moreover, these sub-goals possess semantic characteristics
653 that allow them to be effectively linked with LLM agents.