

Reinforcement Learning for High-Level Strategic Control in Tower Defense Games

Joakim Bergdahl
SEED - Electronic Arts (EA)
Stockholm, Sweden
jbergdahl@ea.com

Alessandro Sestini
SEED - Electronic Arts (EA)
Stockholm, Sweden
asestini@ea.com

Linus Gisslén
SEED - Electronic Arts (EA)
Stockholm, Sweden
lgisslen@ea.com

Abstract—In strategy games, one of the most important aspects of game design is maintaining a sense of challenge for players. Many mobile titles feature quick gameplay loops that allow players to progress steadily, requiring an abundance of levels and puzzles to prevent them from reaching the end too quickly. As with any content creation, testing and validation are essential to ensure engaging gameplay mechanics, enjoyable game assets, and playable levels. In this paper, we propose an automated approach that can be leveraged for gameplay testing and validation that combines traditional scripted methods with reinforcement learning, reaping the benefits of both approaches while adapting to new situations similarly to how a human player would. We test our solution on a popular tower defense game, *Plants vs. Zombies*. The results show that combining a learned approach, such as reinforcement learning, with a scripted AI produces a higher-performing and more robust agent than using only heuristic AI, achieving a 57.12% success rate compared to 47.95% in a set of 40 levels. Moreover, the results demonstrate the difficulty of training a general agent for this type of puzzle-like game.

Index Terms—Reinforcement learning, tower defense, content creation, heuristic AI, game testing

I. INTRODUCTION

The mobile games market is constantly growing and had, as of 2021, the biggest share (50%) of the total gaming market (\$176B) [1]. Modern mobile games can be nearly as content-rich and engaging as AAA computer or console games, which puts new demands on the developers of this type of games that might not have existed a couple of years ago. Many mobile titles are built around quick gameplay loops where the player is progressing through the game at a steady pace. As with any content creation, it needs to be tested and validated. Gameplay mechanics must be engaging, game assets must be enjoyable and levels playable.

This requires an abundance of levels and puzzles to stop the player from reaching the end of the game too quickly. Moreover, many mobile games are live-service titles: long-lasting games that, once published, developers continue to add mechanics, assets, and especially levels. Everything that is added after initial distribution, must be thoroughly tested. One example where this is particularly true is the tower defense genre. Tower defense games are a sub-genre of strategy games, in which the player must defend their territory from waves of

enemies, typically by placing defensive units to counter the attacking forces. Each level presents a puzzle that the player must solve by identifying the right strategies to follow and positioning the appropriate units optimally at each step. Tower defense games are popular in the mobile gaming landscape, with notable examples including *Plants vs. Zombies*, *Bloons TD*, and *Clash Royale*.

One of the big challenges with creating engaging games of this type is the tuning of the level of difficulty. This is especially true for more casual games, such as those in mobile platforms. Often, game designers have just minutes to catch the player's interest and without a well-balanced challenge it could be hard to keep the player from progressing farther into the game. In order for a mobile game to be engaging, it has to have a carefully measured level of difficulty. Too difficult and the gameplay experience becomes frustrating. Too easy and the experience simply becomes unrewarding.

There are many ways of testing the difficulty level of a game. One of the more commonly used methods is to leverage the game's soft-launch phase to gather data on playthroughs and their respective players. However, it is not an ideal way to use live player data to fine-tune the game: it might result in the loss of many potential long-term players that otherwise would have continued to play the game. Another option is to use automated scripts to play the game and in that way estimate the difficulty. A scripted solution is seldom powerful enough to handle new levels and features that are constantly added both in production, and post-production for a live-service games.

Here, we propose a method where we bring together the benefits from human play-testers (adaptive, learning) and the scripted solution (fast, scalable, cheap). In this work, we propose a hybrid approach that combines reinforcement learning (RL) and scripted heuristic AI (HAI). By leveraging the best qualities of the two approaches, one can create a testing AI that is efficient and, at the same time, high-performing, with a higher level of adaptability to changes compared to a HAI-only method. The approach is fairly simple, yet quite powerful: we train the RL agent as a high-level decision-maker, which selects the low-level action the HAI system should perform, deferring the responsibility of how the action is executed to the latter.

We test our approach in a popular and previously mentioned tower defense game: *Plants vs. Zombies*. This game has all the

main characteristics of this genre: a live-service puzzle game played by millions of players, with a simple design but complex gameplay mechanics that require thoughtful strategies. We show how the combination of RL and HAI when trained for each single level performs better than using only the latter. However, we show how this type of puzzle game is hardly generalizable, with each level needing a particular high-level strategy.

II. RELATED WORK

In this section we review the most relevant literature to our contributions.

A. Learning Agents in Games

RL has garnered great interest from the video game community the last few years. Major developments in optimization algorithms and model architectures has lead to impressive results in complex games such as *Star Craft II*, *Dota 2*, and *Gran Turismo 7* [2, 3, 4]. However, these works focused primarily on training the best agent to replace human players, mainly leveraging RL. More recent works explore different methods for training agents in various use cases such as creating believable and human-like game AI [5, 6, 7]. Few approaches have used self-learning agents in mobile games, especially puzzle-like games such as tower defense. Notable examples include: the work by Dias et al. [8], where the authors proposed an RL pipeline for a tower defense player agent; the work by Kristensen et al. [9], which will be examined in a later section; and CandyRL [10], which employs an approach similar to ours for playing Match-3 games.

As mentioned in Section I, automated playtesting through learning agents has gained interest in both research and industry community. However, when dealing with a game in development, training agents with RL seems to be the most feasible solution, particularly during the playtesting phase. A game in development lacks pre-collected datasets, and creating data from expert demonstrations can be time-consuming. Furthermore, it is crucial to be fast and efficient in order to stay in sync with game developers.

B. Automated Playtesting

Many works have employed machine learning for automated gameplay testing [11]. Mugrai et al. [12] showed that mimicked human behavior can be used to achieve more meaningful gameplay testing. Sestini et al. [13] utilized a data-driven technique that allows designers to efficiently train game testing agents, explaining why learning agents are essential for playtesting. Bergdahl et al. [14] proposed a study on the usefulness of using RL policies to playtest levels against scripted agents. Sestini et al. [15] proposed the curiosity-conditioned proximal trajectories algorithm, which tests complex 3D game scenes using a combination of IL, RL, and curiosity-driven exploration. Finally, Abdelfattah et al. [16] proposed an agent that primarily relies on pixel-based state observations while exploring the environment, conditioned on a user’s preference specified by demonstration trajectories.

However, as we will discuss throughout the paper, if one were to employ a RL-only solution for grid-based tower defense games like *Plants vs. Zombies*, the vast number of available actions would make credit assignment particularly difficult. For this reason, we opted for a hybrid approach that combines RL with the pre-existing HAI. Similar approaches have been explored by Karimi et al. [10], Shin et al. [17], and Pang et al. [18], but for different types of games such as Match-3 and real-time strategy games.

C. Difficulty Evaluation

Estimating and evaluating game difficulty is an active topic in video game related research, as it is an important component in good game design. Classically, difficulty has been evaluated and tuned following extensive playtesting by human play-testers, an endeavour that is both time-consuming and expensive. Statistical modeling of difficulty has shown promise in puzzle-style games where the number of moves is limited, successfully predicting the impact on level difficulty when modifying the number of moves allowed [9]. Other researchers have explored using machine learning in a tower defense game to adjust the actual difficulty dynamically in an otherwise static system, managing how enemy waves are spawned [19]. Research applied to the tower defense game *Kingdom Rush: Frontiers* approaches the problem of difficulty assignment through procedural content generation by using flat Monte Carlo search to verify that generated levels are playable, but also that they have a balanced difficulty [20]. All the aforementioned approaches require data collected from well-playing players, that could be difficult to source. As already described in Section I, this paper proposes an approach that combines the benefits from human testers but also from heuristic AI. Developers can leverage on this approach to collect high-quality data for running statistical analyses for difficulty evaluation and validation.

III. ENVIRONMENT

In this section we dive deeper into the *Plants vs. Zombies* (PvZ) game. Figure 1 shows a screenshot of the game.

A. Game Description

In PvZ, the player is tasked to use a set of plants to defend themselves against incoming waves of enemy zombies. What makes PvZ unique is its grid-based nature where the player is limited to place their plants in predefined grids with different layouts of at most five rows, or *lanes*, by nine columns. Similar to other tower defense titles, the complexity of PvZ comes from long-term planning, resource management and the individual abilities of the plant units at the player’s disposal. In each level, PvZ presents all the common tower defense elements: 1) a base that must be defended from attacking enemies; 2) one or more waves of enemies, the diversity and strength of each depending on the particular level and its difficulty; and 3) units with varying defense/offense characteristics required to defend the base against attacking enemies.



Fig. 1. Screenshot of *Plants vs. Zombies 2* gameplay. Enemy zombies attack the player from the right, who has to defend their home base line on the leftmost side of the playarea. If an enemy crosses this line, the player loses. The player is free to place units in any available, unoccupied cell of the game board from the loadout visible on the left. When a unit is greyed out, it is locked behind a cool-down timer. Each unit's sun token cost is displayed in the units lower right corner.

At the start of a round, the player is given a predefined set of plant units, called a *loadout*, they can use in the level. As the level starts, the player is tasked to prepare for the first wave of enemies by placing plant units in unoccupied grid cells. Doing so, the player needs to spend *sun tokens* which are either generated passively in the game-world at a slow rate or through specific plant units that produce them. Each unit incurs a sun token cost, and more capable units tend to be more expensive, demanding planning and prioritization. Importantly, key units will help the player defeat enemies or hinder their advancement. Whenever any unit from the loadout has been placed, this unit will not be usable until a corresponding cool-down timer runs out. This further limits the rate with which players can perform actions, promoting strategy instead. As the enemy waves start spawning, this loop of resource collection and preparation repeats. The termination or game-over criteria for PvZ is fulfilled in two conditions: 1) the player successfully defeats all enemy units and wins; 2) one or more enemy units reach the left-most side of the play field and the player loses.

B. Heuristic AI and Automated Playtesting

When experiments started, the PvZ game had employed an HAI system capable of playing the game to check for stability issues and performance drops. HAI replaces the player and it is composed of two parts: one for computing priority values over a set of strategies, and another for running the strategy with the highest priority. The current implementation of HAI contains 4 of these strategies and we will go into detail about each of them later in Section IV-A. At each step, the HAI takes the current state of the game as input and outputs a priority value for each low-level action. The system then executes the strategy with the higher priority value. We must note that the way the HAI computes the priority values *remains the same across levels*. If the current state is more or less the same but in two different levels, the priority values will be the same. This is a suboptimal solution because, as we will see with our

experiments, some levels require dynamic values to be actually played, and a particular play-style – e.g., preferring defensive plants over attacking ones – can work in one level but not in another. Moreover, if developers want to add a new strategy or a new unit, they need to rebalance all the values.

The HAI system is crucial for PvZ developers, as it allows them to test the game at a very high speed. However, for the reasons listed above, HAI is not general, hardly maintainable, and poorly scalable. Thus, we decided to replace the first system with an RL agent.

IV. METHOD

In this section, we explain our proposed hybrid method that combines HAI and RL, referred to as hybrid RL (HRL) for the remainder of the paper. First, we describe the 4 strategies used by both our hybrid approach and HAI. Then, we describe the algorithm used to train the HRL agent that replaces the priority calculation described in Section III-B and how it is combined with the HAI action execution.

A. High-Level Strategies

The built-in HAI system spans four strategies, capable of covering all gameplay facets of PvZ.

Sow Sun. Central to all tower defense games is the resource management aspect. In PvZ, the main resource collected by the player are sun tokens. Without sun tokens, the agent will not be able to produce or upgrade any units that is needed to complete the level. The *sow sun* strategy is responsible for placing units available in the loadout that contribute to the production of sun tokens on the game board. An example of this is the *sunflower* unit, that spawns sun tokens at fixed intervals.

Attack. In order to successfully complete a game level in tower defense, the player needs to defeat all incoming enemy units, usually split up in individual waves. To do so, the player has a set of offensive units at their disposal. These units vary from close range melee units to ranged shooting units. The *attack* strategy of the classical AI system is responsible for using this combined class of units. In this regard PvZ follows pretty much the typical tower defense setup. An exemplary attack unit is the *peashooter*, that shoots projectiles at approaching enemies in the same lane where it is placed.

Defensive. As resource generation takes time and unit placement incurs a cool-down timer, the player needs to be able to limit enemy advancement. This is done through the *defense* strategy. Units of this type are mainly passive and made to simply block enemies from moving further, with an example being the *walnut*. These units also tend to have a high amount of hit-points to endure damage over time for longer.

Prepare. A special strategy employed by HAI is the *prepare* one. This strategy controls placement of passive units that generally deal no damage nor survives enemy attacks. Their purpose is merely to allow placement of other units of the aforementioned types. An example is the *lilypad* unit which, when placed on a water grid cell, allows for another unit to be placed on-top of it, opening up the game board. Without the lilypad, the water cell does not allow for unit placement.

B. Algorithm

We opted for a combination of RL and HAI, motivated by the following reasons:

- As previously discussed, HAI is difficult to maintain, does not generalize well to new levels and/or mechanics, and has poor scalability. Furthermore, it exhibits suboptimal performance in some levels. Lastly, as we will see in Section VI-C, it lacks robustness when faced with varying levels of difficulty. However, after our initial experiments, we observed that the main issue was primarily due to HAI’s strategy selection system, rather than its execution of actions;
- A fully-fledged RL agent, trained to select one of the units and place it wherever desired, is hardly feasible in this case: each level features a grid of 5×9 cells and a layout of up to 6 units. This results in an action space composed of $5 \times 9 \times 6 = 270$ actions. This presents an efficiency challenge for the RL agent as exploring the action space would be time-consuming, especially when learning something that HAI already manages. In our context, efficiency is crucial, as training new agents should keep pace with game development. Gillberg et al. [21] provide a comprehensive summary of the requirements for applying RL in game development and explain why combining RL with scripted approaches is a good solution in this context.

Hence, we combine the strengths of both approaches while minimizing their respective drawbacks: the hybrid RL agent leverages the knowledge of HAI and determines which low-level strategy the HAI should execute, effectively replacing the strategy selection system of the AI.

In this section, we will go through each of the elements of the employed algorithm. For all the experiments we use proximal policy optimization (PPO) as the optimization method, which is a popular actor-critic on-policy RL algorithm [22]. For this reason, we first have to define two networks: one for the actor and one for the critic.

Model Architecture. For the PPO implementation, the two underlying actor and critic models share no weights and are initialized in separation. Both models are composed of two fully connected layers of size 1024, all with leaky ReLU activation functions. The actor is then followed by an output layer of size 5 fed through a softmax function to produce an action distribution, corresponding to the five actions covering the four available strategies and an additional no-op action. The complete state space is described in Section V-A. More information regarding the action space is provided later in this section. The critic network is followed by an output layer of size 1 with no activation for the value estimate.

Action Masking. Given the resource requirements and cool-down timers restricting gameplay flow, the agent will encounter states where no action will yield an outcome. To avoid these scenarios, action masking is used. At each simulation step, for a given semantic unit type tied to a specific strategy, each unit of said group is checked for. If there is at least

one unit of this group that is cheaper than the available sun token resource pool S and is not currently restricted through a cool-down timer, the action corresponding to this group is labeled as available. This is repeated for each group. In the worst case, no action is available and the game simulation is forwarded to the point where the agent again has agency over the environment, and at least one action is available besides the no-op.

V. EXPERIMENTAL SETUP

In this section we detail the experimental setup used to evaluate our hybrid approach. All experiments utilizing 5 different seeds were run over 5 machines, each with an Intel Xeon Gold 6230 @ 2.10 GHz CPU, 128 GB of RAM and an NVIDIA RTX A6000 GPU with 11 GB of VRAM. On average, a training run for one model completed in 5 hours and each model was trained for 10K episodes.

A. Environment Details

In this section we detail the following components: observation space, action space, and reward function.

Observation Space. This information is represented as feature vector observations of size 88, composed by:

- One-hot encodings of the semantic type of each unit in the loadout with masking if the loadout is incomplete. The maximum number of units in the loadout is 6, and some levels can have less than 6 units. In total, for this we have 6×6 values;
- Current cool-down statuses of each loadout unit. Each value is in $[0, 1]$ with 0 meaning the unit is available, and for this we have a total of 6 values;
- Amount of currently held sun tokens, only 1 value;
- Average health points over all enemy units for each row, for a total of 5 values;
- The distances to the closest enemy in each row, for a total of 5 values;
- Number of enemies that have advanced past the mid-point of the board per row, for a total of 5 values;
- Number of enemies per row, for a total of 5 values; and
- Count of planted units of each semantic type for each row, for a total of 5×5 values.

Action Space and Execution. Deferring low-level strategy execution to the classical AI system, the action space for the HRL agent is simple. The agent only needs to produce one out of four available high-level strategies at each decision step as described in Section IV-A. In addition, a no-op action is introduced to allow for accumulating resources or waiting for more optimal actions. As the game might be in a state where no action can be performed due to lack of sun tokens or all plants being unavailable due to cool-downs, the simulation is progressed until at least one additional action than the no-op is available through the action mask as detailed in Section IV-B.

Reward Function. Tower defense games require the player to survive and defeat waves of incoming enemies. As such, the agent gets a positive reward of $1.2 \cdot N$ where N is the

TABLE I

PERFORMANCE COMPARISON OF SUCCESS RATES, EPISODIC REWARDS AND EPISODE LENGTHS IN STEPS BETWEEN HRL, HAI AND A RANDOM AGENT. ALL VALUES REPRESENT THE MEAN OF 100 EPISODES WITH 5 SEEDS. FOR HRL, WE TRAIN ONE AGENT FOR EACH LEVEL. HRL OUTPERFORMS THE BASELINES IN MOST OF THE LEVELS, ACHIEVING BOTH HIGHER SUCCESS RATES AND REWARDS. IN THE FINAL ROW, THE AVERAGE OF THE SUCCESS RATE AND THE SUM OF THE REWARDS AND STEPS FOR EACH LEVEL ARE REPORTED. NOTE, THE LEVELS ARE NOT ORDERED BY COMPLEXITY AND THERE IS ONE PARTICULAR LEVEL (LEVEL 23) WHERE HAI ACHIEVES A VERY LOW REWARD.

Level	Success Rate (%)						Cumulative Reward (R)						Average Number of Steps					
	HAI		Random		HRL		HAI		Random		HRL		HAI		Random		HRL	
1	14.00 ± 3.90	6.60 ± 3.01	29.80 ± 15.66	11.56 ± 0.82	0.63 ± 0.93	11.40 ± 5.23	351.02 ± 9.40	263.88 ± 6.92	404.37 ± 64.51									
2	45.80 ± 2.40	5.40 ± 0.80	67.00 ± 9.61	8.41 ± 0.48	-5.96 ± 0.21	10.72 ± 1.49	246.57 ± 3.52	143.71 ± 3.40	266.55 ± 19.17									
3	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	-7.53 ± 0.34	-7.66 ± 0.31	5.93 ± 0.67	410.37 ± 8.33	388.69 ± 7.16	376.99 ± 50.17								
4	0.00 ± 0.00	0.00 ± 0.00	3.40 ± 1.74	2.41 ± 0.25	-0.54 ± 0.56	5.24 ± 1.38	349.45 ± 2.59	309.37 ± 8.23	425.46 ± 54.24									
5	8.00 ± 1.67	3.00 ± 1.10	17.20 ± 3.76	3.79 ± 0.44	-5.22 ± 0.43	7.85 ± 0.43	256.55 ± 7.08	180.83 ± 3.25	275.75 ± 21.82									
6	85.40 ± 1.36	49.00 ± 2.19	86.60 ± 11.15	15.11 ± 0.40	-3.13 ± 0.86	15.45 ± 7.25	422.45 ± 2.88	533.93 ± 7.24	446.33 ± 41.71									
7	100.00 ± 0.00	94.00 ± 1.67	99.20 ± 0.40	17.76 ± 0.10	13.94 ± 0.29	17.03 ± 0.40	211.61 ± 2.31	165.69 ± 31.39	151.46 ± 17.29									
8	99.20 ± 0.75	96.60 ± 1.85	100.00 ± 0.00	15.38 ± 0.15	13.13 ± 0.30	16.93 ± 0.13	245.62 ± 2.67	619.04 ± 101.28	280.48 ± 26.68									
9	28.40 ± 4.32	41.20 ± 5.74	40.20 ± 13.48	7.69 ± 0.88	8.19 ± 1.14	9.10 ± 2.01	271.52 ± 5.72	283.35 ± 6.77	282.38 ± 27.85									
10	99.60 ± 0.49	100.00 ± 0.00	99.80 ± 0.40	21.77 ± 0.14	21.80 ± 0.18	22.04 ± 0.26	415.36 ± 1.08	423.91 ± 1.92	432.09 ± 16.59									
11	94.80 ± 2.04	98.60 ± 0.80	100.00 ± 0.00	12.95 ± 0.18	12.27 ± 0.40	15.63 ± 0.18	305.55 ± 5.29	409.92 ± 53.77	231.44 ± 32.11									
12	38.20 ± 4.96	12.80 ± 2.48	46.00 ± 9.38	5.01 ± 0.61	-3.60 ± 0.67	6.17 ± 1.54	231.40 ± 4.78	194.54 ± 10.51	251.01 ± 26.93									
13	0.00 ± 0.00	0.20 ± 0.40	3.40 ± 3.77	3.61 ± 0.36	-5.71 ± 0.65	4.16 ± 4.20	282.41 ± 4.42	181.45 ± 5.03	329.89 ± 44.73									
14	35.00 ± 3.03	55.80 ± 3.43	59.40 ± 12.03	11.42 ± 0.33	14.13 ± 0.89	15.70 ± 2.14	364.00 ± 6.54	466.59 ± 66.48	453.01 ± 32.41									
15	41.00 ± 3.22	27.60 ± 2.42	67.00 ± 13.07	20.19 ± 1.30	12.90 ± 0.72	26.48 ± 4.44	342.71 ± 13.62	330.68 ± 1.74	393.13 ± 63.84									
16	41.00 ± 3.35	2.40 ± 0.49	53.40 ± 14.11	24.21 ± 0.67	5.18 ± 0.97	28.12 ± 3.38	465.32 ± 7.87	422.91 ± 9.81	499.67 ± 24.40									
17	39.60 ± 7.68	8.60 ± 2.15	47.00 ± 11.01	0.78 ± 1.53	-10.41 ± 0.80	4.38 ± 4.41	325.33 ± 7.47	269.51 ± 26.27	349.38 ± 35.32									
18	100.00 ± 0.00	100.00 ± 0.00	99.80 ± 0.40	4.36 ± 0.00	5.31 ± 0.04	5.82 ± 0.06	73.00 ± 0.00	216.51 ± 123.01	48.51 ± 9.84									
19	100.00 ± 0.00	71.00 ± 6.45	99.60 ± 0.80	6.39 ± 0.00	2.80 ± 0.91	9.94 ± 0.93	207.00 ± 0.00	235.80 ± 10.20	154.23 ± 16.80									
20	100.00 ± 0.00	100.00 ± 0.00	100.00 ± 0.00	-74.62 ± 0.00	-6.94 ± 0.32	4.38 ± 0.12	85.00 ± 0.00	45.30 ± 0.63	49.35 ± 7.46									
21	90.00 ± 0.89	97.00 ± 1.10	97.60 ± 1.36	17.08 ± 0.35	19.25 ± 0.50	21.22 ± 0.25	415.73 ± 5.79	583.79 ± 9.01	388.65 ± 29.82									
22	0.00 ± 0.00	0.00 ± 0.00	13.80 ± 10.81	-2.96 ± 0.34	-5.21 ± 0.51	1.99 ± 3.10	303.74 ± 2.09	124.65 ± 3.06	353.64 ± 51.16									
23	0.00 ± 0.00	0.00 ± 0.00	5.60 ± 11.20	-174.88 ± 0.00	-32.07 ± 0.85	-2.58 ± 7.38	255.00 ± 0.00	77.38 ± 2.93	273.62 ± 106.08									
24	68.60 ± 2.94	56.60 ± 11.07	73.40 ± 12.82	19.25 ± 0.71	13.73 ± 2.26	18.67 ± 3.61	338.93 ± 6.20	162.99 ± 5.98	194.64 ± 66.75									
25	100.00 ± 0.00	100.00 ± 0.00	100.00 ± 0.00	10.25 ± 0.00	10.90 ± 0.03	11.36 ± 0.03	89.00 ± 0.00	48.27 ± 12.59	31.92 ± 2.84									
26	0.00 ± 0.00	6.60 ± 1.96	100.00 ± 0.00	-6.87 ± 0.00	-7.49 ± 0.45	10.08 ± 0.08	121.00 ± 0.00	70.75 ± 11.96	57.97 ± 11.15									
27	69.40 ± 1.50	63.00 ± 2.10	91.40 ± 5.61	15.53 ± 0.23	13.00 ± 0.33	21.09 ± 1.05	325.88 ± 2.04	467.96 ± 12.52	406.18 ± 23.73									
28	100.00 ± 0.00	93.40 ± 1.36	99.60 ± 0.49	19.64 ± 0.09	16.01 ± 0.30	18.60 ± 0.61	131.32 ± 1.61	195.00 ± 26.80	132.58 ± 9.55									
29	2.20 ± 1.94	1.20 ± 1.47	5.60 ± 1.85	-0.01 ± 0.43	-5.83 ± 0.70	1.28 ± 0.90	228.97 ± 7.19	176.13 ± 5.13	220.15 ± 9.85									
30	76.40 ± 5.35	73.80 ± 2.71	98.00 ± 2.28	17.91 ± 0.70	15.35 ± 0.80	25.13 ± 0.98	276.06 ± 5.83	397.40 ± 29.53	327.08 ± 56.40									
31	100.00 ± 0.00	70.60 ± 3.38	99.00 ± 0.63	24.87 ± 0.40	17.67 ± 0.61	24.58 ± 1.41	338.47 ± 4.62	516.83 ± 16.29	352.61 ± 52.40									
32	1.20 ± 1.60	2.60 ± 2.58	6.80 ± 1.47	-3.62 ± 1.26	-5.81 ± 1.52	-0.51 ± 2.09	422.42 ± 10.68	362.50 ± 17.40	447.01 ± 60.89									
33	3.40 ± 1.85	0.60 ± 0.80	18.40 ± 15.87	4.55 ± 0.85	-5.10 ± 1.02	13.04 ± 6.29	278.93 ± 6.58	188.86 ± 5.56	370.10 ± 77.41									
34	100.00 ± 0.00	97.00 ± 1.41	99.80 ± 0.40	19.15 ± 0.10	16.90 ± 0.08	19.44 ± 0.78	258.76 ± 1.52	371.30 ± 14.01	264.62 ± 27.82									
35	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	-7.82 ± 0.31	-15.14 ± 0.53	-0.87 ± 0.93	233.28 ± 4.51	70.94 ± 0.81	370.25 ± 50.61									
36	58.00 ± 5.14	36.40 ± 4.08	91.20 ± 3.66	12.98 ± 0.74	2.44 ± 1.34	18.55 ± 0.44	358.09 ± 5.21	419.10 ± 16.41	427.86 ± 19.32									
37	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	-13.98 ± 0.25	-15.42 ± 0.54	-8.94 ± 0.14	442.27 ± 2.88	450.25 ± 9.16	417.11 ± 8.19									
38	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	-4.00 ± 0.00	-27.94 ± 0.43	-11.37 ± 0.01	119.00 ± 0.00	180.31 ± 3.84	308.55 ± 21.94									
39	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00	-10.61 ± 0.14	-13.00 ± 0.18	-7.51 ± 0.86	172.85 ± 1.42	155.41 ± 1.49	292.31 ± 54.21									
40	79.00 ± 0.89	0.00 ± 0.00	65.60 ± 18.38	17.77 ± 0.30	-20.78 ± 0.90	16.06 ± 6.11	514.86 ± 3.07	337.60 ± 4.09	480.13 ± 14.61									
Total:	47.95 ± 1.53	39.29 ± 1.73	57.12 ± 5.19	64.88 ± 16.18	32.57 ± 25.46	419.92 ± 77.70	11486.80 ± 166.81	11443.03 ± 693.58	12218.46 ± 1388.60									

number of enemies defeated since the last action decision. To prolong play, allowing for more complicated action selection, the agent is also penalized with a fixed value of $r_Z = -1/200$ for each enemy that has been allowed to advance since the last action decision. The reasoning behind this is to make use of the defensive strategy more appealing, leading to more placement of blocking plant units. Finally, the agent receives a terminal reward of either 1 or -5 depending on if the level was successfully completed or the agent lost.

B. Levels and Difficulty

For our experiments, we decided to use a set of 40 already existing levels. Although they do not encompass all the levels in the game, this set summarizes all the gameplay mechanics of PvZ. These levels span across very easy and more difficult levels, with varied gameplay mechanics and various types and quantities of loadouts. For our generalization experiments, we train our agent in subsets of the original set of 40 levels. We train the generalization agents with different amounts of levels: 5, 10, and 20. These subsets are randomly sampled from the original set. As we will see in Section VI-B, we then test the generalization agents in the full set of 40 levels.

As mentioned in Section I, PvZ, as well as most tower defense games, includes a difficulty system. This system is used to keep the player engaged with the game: for more advanced players, the level should be more difficult than for

less experienced ones. For this reason, in each level, we can set a level of difficulty. For our experiments, we train all of our agents for all levels with a fixed difficulty of 100K, which is generally high and makes most of the levels beatable only using mindful strategies. The value of the difficulty ultimately affects the composition, strength and health of the enemies in each wave. As we will see in Section VI-C, to test the robustness of our agent, we vary the level of difficulty during testing between 0 and 200K, with the latter representing an almost impossible level to complete.

C. Experiments and Baselines

For this evaluation, we are mainly interested in 4 research questions:

- Is our proposed approach able to solve more levels than baselines?
- Is our proposed approach more robust to increasing game difficulty?
- Can our proposed approach generalize to unseen levels during training and still outperform baselines?

We will answer these questions later in Section VI. For comparison, two baselines were considered.

- **Random:** we replace the priority calculator in HAI with a random number generator, sampling strategies following action masking as mentioned in Sections IV-B;

TABLE II

PERFORMANCE OF AGENTS TRAINED ON SUBSETS OF THE MAIN LEVELSET OF 40 LEVELS VS. HAI. THE AGENTS WERE TRAINED ON N RANDOMLY SAMPLED LEVELS WHERE $N \in \{5, 10, 20\}$. VALUES IN THE HRL[†] COLUMN ARE IDENTICAL TO THE ONES IN TABLE I, WHERE ONE MODEL WAS TRAINED FOR EACH LEVEL RESPECTIVELY. THE FINAL ROW PRESENTS THE MEAN OVER ALL SUCCESS RATES AND THE SUM OF THE EPISODIC REWARDS.

Level	HAI		HRL [†]		HRL ($N = 5$)		HRL ($N = 10$)		HRL ($N = 20$)	
	Success Rate	Reward	Success Rate	Reward	Success Rate	Reward	Success Rate	Reward	Success Rate	Reward
1	14.00 ± 3.90	11.56 ± 0.82	29.80 ± 15.66	11.40 ± 5.23	7.60 ± 2.65	7.93 ± 2.26	14.60 ± 7.55	6.54 ± 5.50	6.60 ± 4.59	7.17 ± 2.73
2	45.80 ± 2.40	8.41 ± 0.48	67.00 ± 9.61	10.72 ± 1.49	56.40 ± 10.48	10.22 ± 1.47	43.60 ± 24.98	6.05 ± 6.85	32.20 ± 15.45	6.62 ± 2.32
3	0.00 ± 0.00	-7.53 ± 0.34	0.00 ± 0.00	-5.93 ± 0.67	0.00 ± 0.00	-5.50 ± 0.42	0.00 ± 0.00	-6.87 ± 0.64	0.00 ± 0.00	-5.54 ± 0.66
4	0.00 ± 0.00	2.41 ± 0.25	3.40 ± 1.74	5.24 ± 1.38	1.00 ± 0.63	1.79 ± 1.21	1.00 ± 1.55	-0.54 ± 2.19	0.00 ± 0.00	0.38 ± 1.82
5	8.00 ± 1.67	3.79 ± 0.44	17.20 ± 3.76	7.85 ± 0.43	9.80 ± 2.99	5.91 ± 1.24	13.00 ± 6.32	3.70 ± 4.89	5.80 ± 3.54	4.48 ± 1.70
6	85.40 ± 1.36	15.11 ± 0.40	86.60 ± 11.15	15.45 ± 7.25	75.60 ± 17.85	8.96 ± 8.33	56.40 ± 11.94	1.19 ± 4.24	88.60 ± 6.62	15.16 ± 2.99
7	100.00 ± 0.00	17.76 ± 0.10	99.20 ± 0.40	17.03 ± 0.40	97.60 ± 2.50	16.56 ± 0.83	95.00 ± 3.29	14.60 ± 1.09	97.40 ± 1.36	15.87 ± 1.30
8	99.20 ± 0.75	15.38 ± 0.15	100.00 ± 0.00	16.93 ± 0.13	99.00 ± 1.55	14.45 ± 1.06	96.20 ± 4.96	13.10 ± 2.26	89.60 ± 11.22	13.51 ± 2.06
9	28.40 ± 4.32	7.69 ± 0.88	40.20 ± 13.48	9.10 ± 2.01	20.00 ± 19.73	3.55 ± 5.04	38.80 ± 12.20	8.04 ± 2.28	33.80 ± 17.85	5.87 ± 5.53
10	99.60 ± 0.49	21.77 ± 0.14	99.80 ± 0.40	22.04 ± 0.26	99.80 ± 0.40	21.60 ± 0.09	79.60 ± 39.80	15.72 ± 11.86	96.20 ± 7.11	20.76 ± 1.79
11	94.80 ± 2.04	12.95 ± 0.18	100.00 ± 0.00	15.63 ± 0.18	99.20 ± 0.75	13.71 ± 0.93	95.60 ± 6.37	11.95 ± 1.82	97.20 ± 2.79	12.86 ± 1.15
12	38.20 ± 4.96	5.01 ± 0.61	46.00 ± 9.38	6.17 ± 1.54	35.00 ± 6.48	4.92 ± 1.12	36.00 ± 15.23	3.28 ± 4.16	32.80 ± 11.46	4.72 ± 1.75
13	0.00 ± 0.00	3.61 ± 0.36	3.40 ± 3.77	4.16 ± 4.20	0.60 ± 0.49	3.03 ± 1.01	0.40 ± 0.49	-0.01 ± 3.46	0.00 ± 0.00	2.42 ± 1.61
14	35.00 ± 3.03	11.42 ± 0.33	59.40 ± 12.03	15.70 ± 2.14	14.00 ± 2.53	8.12 ± 0.38	31.20 ± 19.88	9.54 ± 3.96	17.40 ± 2.42	8.36 ± 0.63
15	41.00 ± 3.22	20.19 ± 1.30	67.00 ± 13.07	26.48 ± 4.44	20.60 ± 16.17	9.44 ± 7.58	33.60 ± 13.17	15.90 ± 3.69	30.00 ± 15.66	12.72 ± 8.23
16	41.00 ± 3.35	24.21 ± 0.67	53.40 ± 14.11	28.12 ± 3.38	15.60 ± 12.63	13.74 ± 6.99	4.20 ± 2.32	7.79 ± 3.10	33.00 ± 12.08	22.69 ± 4.56
17	39.60 ± 7.68	0.78 ± 1.53	47.00 ± 11.01	4.38 ± 4.41	32.80 ± 5.27	1.85 ± 2.51	27.20 ± 11.32	-3.25 ± 4.77	27.20 ± 12.56	2.34 ± 2.26
18	100.00 ± 0.00	4.36 ± 0.00	99.80 ± 0.40	5.82 ± 0.06	99.60 ± 0.49	4.97 ± 0.56	96.60 ± 6.80	4.83 ± 0.55	89.60 ± 16.67	4.84 ± 0.49
19	100.00 ± 0.00	6.39 ± 0.00	99.60 ± 0.80	9.94 ± 0.93	77.40 ± 24.93	6.94 ± 4.03	99.20 ± 0.98	9.60 ± 1.59	98.00 ± 4.00	9.90 ± 1.54
20	100.00 ± 0.00	-74.62 ± 0.00	100.00 ± 0.00	4.38 ± 0.12	100.00 ± 0.00	-5.55 ± 3.67	100.00 ± 0.00	-1.85 ± 5.13	100.00 ± 0.00	2.48 ± 6.05
21	90.00 ± 0.89	17.08 ± 0.35	97.60 ± 1.36	21.22 ± 0.25	66.60 ± 20.71	12.14 ± 3.96	69.40 ± 14.68	12.33 ± 3.37	64.40 ± 20.04	11.37 ± 3.97
22	0.00 ± 0.00	-2.96 ± 0.34	13.80 ± 10.81	1.99 ± 3.10	0.20 ± 0.40	-3.51 ± 0.71	1.80 ± 3.60	-3.96 ± 0.98	0.00 ± 0.00	-4.17 ± 0.36
23	0.00 ± 0.00	-174.88 ± 0.00	5.60 ± 11.20	-2.58 ± 7.38	0.20 ± 0.40	-37.26 ± 19.34	1.20 ± 1.94	-24.38 ± 14.45	0.00 ± 0.00	-18.93 ± 20.62
24	68.60 ± 2.94	19.25 ± 0.71	73.40 ± 12.82	18.67 ± 3.61	70.40 ± 8.78	17.97 ± 2.25	58.00 ± 14.49	14.01 ± 4.97	59.20 ± 16.03	14.91 ± 5.10
25	100.00 ± 0.00	10.25 ± 0.00	100.00 ± 0.00	11.36 ± 0.03	100.00 ± 0.00	10.69 ± 0.43	99.40 ± 1.20	10.61 ± 0.50	100.00 ± 0.00	10.75 ± 0.49
26	0.00 ± 0.00	-6.87 ± 0.00	100.00 ± 0.00	10.08 ± 0.08	74.20 ± 35.77	4.83 ± 7.25	45.20 ± 45.11	-0.56 ± 8.81	45.40 ± 42.15	-1.03 ± 8.39
27	69.40 ± 1.50	15.53 ± 0.23	91.40 ± 5.61	21.09 ± 1.05	74.20 ± 8.13	16.86 ± 1.61	66.80 ± 20.74	14.54 ± 6.15	52.40 ± 17.72	12.93 ± 3.63
28	100.00 ± 0.00	19.64 ± 0.09	99.60 ± 0.49	18.60 ± 0.61	99.20 ± 1.17	18.61 ± 0.49	98.60 ± 2.80	18.13 ± 1.47	95.20 ± 6.73	17.72 ± 1.61
29	2.20 ± 1.94	-0.01 ± 0.43	5.60 ± 1.85	1.28 ± 0.90	1.60 ± 1.62	0.37 ± 1.15	2.80 ± 2.04	-1.33 ± 3.11	0.80 ± 0.98	-0.41 ± 1.05
30	76.40 ± 5.35	17.91 ± 0.70	98.00 ± 2.28	25.13 ± 0.98	51.60 ± 8.85	14.34 ± 1.41	61.40 ± 14.50	15.00 ± 1.86	39.40 ± 8.82	12.57 ± 1.60
31	100.00 ± 0.00	24.87 ± 0.40	99.00 ± 0.63	24.58 ± 1.41	99.60 ± 0.80	24.02 ± 0.70	90.20 ± 10.96	22.24 ± 3.34	99.60 ± 0.49	24.55 ± 0.50
32	1.20 ± 1.60	-3.62 ± 1.26	6.80 ± 1.47	-0.51 ± 2.09	1.00 ± 1.10	-5.79 ± 1.61	0.80 ± 0.75	-4.73 ± 1.86	2.80 ± 2.40	-3.89 ± 3.23
33	3.40 ± 1.85	4.55 ± 0.85	18.40 ± 15.87	13.04 ± 6.29	3.00 ± 1.10	4.94 ± 0.94	4.40 ± 2.50	3.31 ± 4.53	0.40 ± 0.49	2.99 ± 1.74
34	100.00 ± 0.00	19.15 ± 0.10	99.80 ± 0.40	19.44 ± 0.78	100.00 ± 0.00	16.81 ± 2.00	99.80 ± 0.40	17.00 ± 1.68	100.00 ± 0.00	18.29 ± 0.38
35	0.00 ± 0.00	-7.82 ± 0.31	0.00 ± 0.00	-0.87 ± 0.93	0.00 ± 0.00	-8.93 ± 4.57	0.00 ± 0.00	-12.48 ± 5.27	0.00 ± 0.00	-12.29 ± 5.17
36	58.00 ± 5.14	12.98 ± 0.74	91.20 ± 3.66	18.55 ± 0.44	38.80 ± 13.53	9.66 ± 2.90	50.80 ± 17.01	9.42 ± 4.79	38.40 ± 19.58	10.54 ± 2.89
37	0.00 ± 0.00	-13.98 ± 0.25	0.00 ± 0.00	-8.94 ± 0.14	0.00 ± 0.00	-8.90 ± 0.32	0.00 ± 0.00	-11.06 ± 2.37	0.00 ± 0.00	-9.04 ± 0.19
38	0.00 ± 0.00	-4.00 ± 0.00	0.00 ± 0.00	-11.37 ± 0.01	0.00 ± 0.00	-100.17 ± 112.53	0.00 ± 0.00	-56.75 ± 89.70	0.00 ± 0.00	-11.56 ± 0.28
39	0.00 ± 0.00	-10.61 ± 0.14	0.00 ± 0.00	-7.51 ± 0.86	0.00 ± 0.00	-9.06 ± 0.19	0.00 ± 0.00	-10.16 ± 1.32	0.00 ± 0.00	-9.01 ± 0.12
40	79.00 ± 0.89	17.77 ± 0.30	65.60 ± 18.38	16.06 ± 6.11	15.60 ± 11.11	-3.05 ± 7.16	3.40 ± 4.22	-13.72 ± 5.90	15.60 ± 20.58	-6.61 ± 12.70
Total:	47.95 ± 1.53	64.88 ± 16.18	57.12 ± 5.19	419.92 ± 77.70	43.94 ± 6.05	121.21 ± 222.25	42.90 ± 8.65	116.77 ± 240.46	42.22 ± 7.53	227.29 ± 125.19

- **HAI**: the full heuristic AI that consists of a hand-crafted priority calculator and low-level strategy executors. More details in Section III-B.

VI. RESULTS

In this section, experimental results are reported. Each experiment was repeated using different seeds for 5 times over which the reported results are averaged and presented with mean and standard deviation. In all the figures and tables, our hybrid approach is identified as *HRL*.

A. General Performance

In the first experiment, we train one HRL agent for each of the 40 levels in the training suite and compared these results against HAI and a random agent. Note that these levels are not in order of complexity, meaning that level 40 is not necessarily more complex than level 1. The results are reported in Table I. Each of these agents was trained for 10K episodes at a fixed level of difficulty of 100K – that is generally high for each level. As the table indicates, our approach achieves a generally higher success rate compared to the random agent and especially HAI. For the reasons detailed in Section III-B, the static strategy selection system is not optimal for playing all levels in this game. As the table shows, some levels are quite challenging to complete (e.g., levels 37, 38, and 39), with none of the approaches being able to solve them. However, in these levels our hybrid approach attains a higher reward compared to the others, meaning that it can eliminate more zombies and survive longer. Interestingly, there are levels (e.g.,

levels 10, 18, and 25) where even a random agent can achieve a 100% success rate. In some instances (e.g., level 18), our agent attains a slightly lower success rate (99.80%). Notably, there are some levels (e.g., levels 9, 21, and 26) where HAI performs significantly worse than a random agent, while this never occurs for HRL. Nonetheless, even though our hybrid approach achieves a higher overall success rate, we observe that HAI slightly outperforms HRL in some levels (e.g. level 7, 31, and 34). In total, our hybrid approach achieves a success rate of 57.12%, compared to 47.95% of HAI and 39.29% of a random agent. At the same time, our agent achieves a much higher reward compared to the baselines: a total of 419.92 compared to 64.88 of HAI and 32.57 of the random agent. Moreover, our agent has a higher total number of timesteps, meaning that it survives longer than the baselines. Even if there is not always a correlation between the number of steps and success rate, this value combined with the reward indicates that our agent survives longer *and* kills more zombies, two important characteristics to win the game. In case of HAI, there is one particular level (level 23) where this baseline achieves a very low reward that decreases the total value. However, even if we remove this outlier, our approach still achieves a higher reward.

In Figure 2, we show the training progression of the agent across four levels: levels 2, 6, 15, and 24. We compare the training progression with the success rates of HAI and the random agent. As the baselines are not trained, their values remain constant regardless of the training. From these plots we can see that our agent starts to outperform HAI in

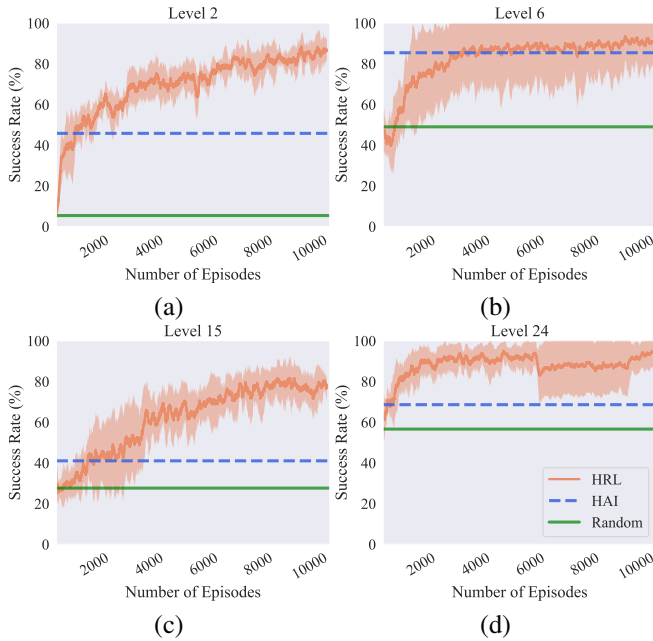


Fig. 2. Training progression of the HRL agent (orange) over four levels, compared to the mean success rate of HAI (blue) and a random agent (green) over 100 episodes in each level respectively. In the displayed levels, the HRL agent learns to outperform HAI in approximately 400 to 1500 episodes. Both the HRL agent and HAI consistently outperform the random agent. The gathered statistics are averaged over 5 different seeds.

approximately 400 to 1500 training episodes.

It is interesting to compare the general behavior between HAI and our hybrid approach. Figure 3 shows the action distributions for four example levels: levels 2, 6, 15, and 24. As illustrated by the plots, HAI exhibits a considerably more aggressive behavior in all the levels, with no significant differences between behaviors across different levels. In contrast, our hybrid approach demonstrates a more cautious behavior in all levels, but in varying ways: in levels 2 and 15, our hybrid approach waits for the opportune moment to deploy a defensive unit, while in levels 6 and 24, it plants additional sunflower units to gather more resources for deploying more effective offensive units.

B. Generalization Performance

In this experiment, we train a set of three agents: one trained in 5 levels, another in 10 levels, and the last one in 20 levels. All agents were trained for 10K episodes. The levels chosen for training these agents were randomly sampled from the original list of 40 levels. We repeat this experiment for 5 seeds, each time sampling a different subset of levels. In contrast to Section VI-A, where we trained one agent for each level, these agents were exposed to multiple levels during training and then tested on unseen levels without retraining. Our goal is to evaluate the ability of our approach to generalize to unseen levels and determine the number of levels needed for training an agent that generalizes effectively. Table II presents the results of this experiment. As the table indicates, these agents are unable to generalize sufficiently to outperform the agents

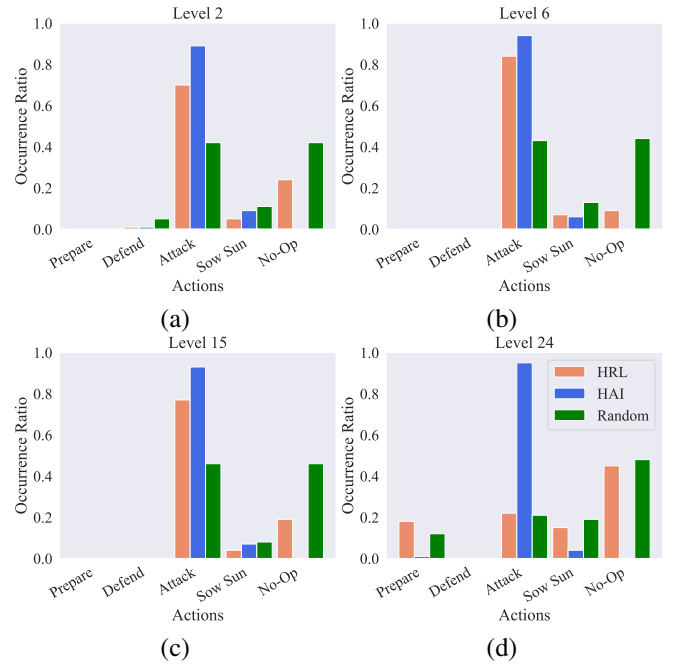


Fig. 3. Action distribution comparison between HAI, random and HRL agents over four levels representing the average, normalized occurrence of each action over 100 episodes repeated for 5 different seeds. As evident from the plots, the HRL agent utilizes the actions differently than the other baselines from level to level, meaning the agent has learned to better leverage actions based on the dynamics of the level.

trained in Section VI-A, and especially HAI. This could imply that each level requires a specific, non-transferable strategy. This observation reflects one of the main characteristics of tower defense games: each level is a puzzle that necessitates its own strategy, and players must fail before discovering the correct one. This is precisely what our agents trained in Section VI-A do: they fail and learn for each level.

C. Difficulty Robustness

In this experiment, we use the agents trained in Section VI-A. The goal of this experiment is to determine whether the performance of trained agents changes when the level of difficulty of a particular puzzle is altered, and how it compares to HAI. We conduct this experiment for four different levels – level 2, 6, 15, and 24 – and for five different difficulty levels – [0, 50K, 100K, 150K, 200K]. We have chosen these levels as they are levels where all the baselines perform reasonably well at difficulty level 100K (see Table I). Figure 4 displays the results. As evident from the figure, our hybrid RL solution is capable of maintaining higher performance even at higher difficulties compared to HAI. This demonstrates that our agent finds a better solution for a specific puzzle regardless of its difficulty level, in contrast to HAI.

VII. CONCLUSION AND DISCUSSION

In this paper, we propose a hybrid approach that combines RL and heuristic AI, thus integrating the performance of an RL-trained agent with the scalability of hard-coded scripted AI. We tested our approach in the popular tower defense game,

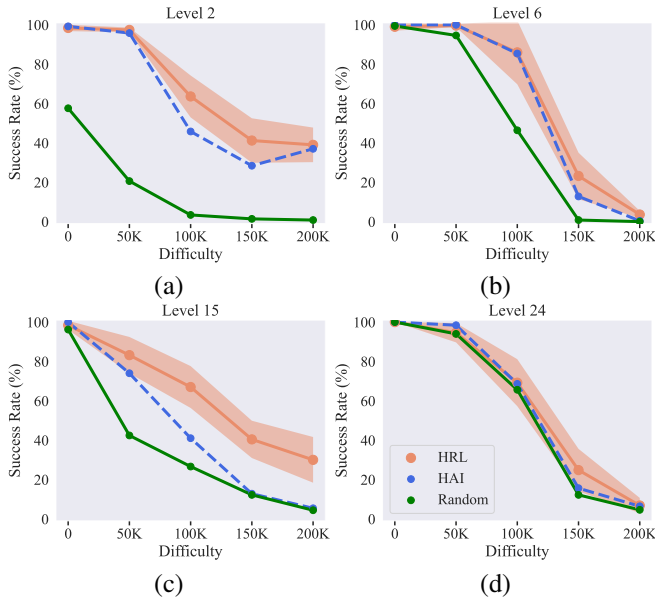


Fig. 4. Performance of HRL agent (orange) compared to HAI (blue) and a random agent (green) with increasing difficulty in 4 levels. Five difficulties were used, specifically [0, 50K, 100K, 150K, 200K]. More details about the difficulty can be found in Section V-B. The general performance of the agent indicates that it successfully completes more levels than the baselines for each difficulty. We collect statistics for the experiments for 100 levels and 5 different seeds. Note, in the experiments, the agents were trained with difficulty of 100K, the mid-point of the aforementioned difficulty set.

Plants vs. Zombies. We found that leveraging existing expertise and knowledge of a classical AI system and augmenting its capabilities using a self-learning strategy prioritizer leads to improved performance over the original AI system.

Our experiments show that this approach can approximate high-performing players, but also their adaptability, that could be leveraged to generate data useful for level testing and validation at scale. With little human intervention, we show that the approach can be applied to previously unseen levels simply through re-training. However, we also conclude that the nature of tower defense games, which are very puzzle-like in their nature, makes it hard to train a general solution for all levels. The recommendation is to rather train a model for each level to get consistent results.

Although our approach demonstrated promising results in the tested game, we believe there is still room for improvement. As we saw little evidence of generalization over previously unseen levels both for the HAI and the hybrid solution important future work here would be to create an algorithm that can handle levels/features/situations that are fundamentally different without the need of retraining. We also see that other game genres, which uses a HAI for low-level control but also have a strategic component that can be replaced by RL, can benefit from this proposed solution. For instance, we believe real-time strategy games would be of interest for future application of this approach. Even for team based action games with a capable low-level heuristic AI system for locomotion, navigation, and combat, this approach could be used to improve high-level strategy by choosing

which low-level behavioral routine the bots should run.

ACKNOWLEDGMENTS

We would like to express our gratitude to Timur Solovev and Mónica Villanueva Aylagas for their valuable contributions, as well as Cory Deines, Rain Zhang and Geoff Schemmel from the *PopCap* team for their extensive support.

REFERENCES

- [1] D. M. Bernes, "A comparison of product placement in video games: aesthetics matter more than value." Ph.D. dissertation, University of Applied Sciences, 2022.
- [2] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik *et al.*, "Grandmaster level in StarCraft II using multi-agent reinforcement learning," *Nature*, 2019.
- [3] C. Berner, G. Brockman, B. Chan, V. Cheung, P. Debiak, C. Dennison, D. Farhi *et al.*, "Dota 2 with large scale deep reinforcement learning," *arXiv preprint 1912.06680*, 2019.
- [4] P. R. Wurman, S. Barrett, K. Kawamoto, J. MacGlashan, K. Subramanian, T. J. Walsh, R. Capobianco *et al.*, "Outracing champion gran turismo drivers with deep reinforcement learning," *Nature*, 2022.
- [5] S. Milani, A. Juliani, I. Momennejad, R. Georgescu, J. Rzepecki, A. Shaw, G. Costello, F. Fang, S. Devlin, and K. Hofmann, "Navigates like me: Understanding how people evaluate human-like ai in video games," in *Conference on Human Factors in Computing Systems*, 2023.
- [6] T. Pearce and J. Zhu, "Counter-strike deathmatch with large-scale behavioural cloning," in *Conference on Games*, 2022.
- [7] G. Wang, Y. Xie, Y. Jiang, A. Mandelkar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar, "Voyager: An open-ended embodied agent with large language models," *arXiv preprint arXiv:2305.16291*, 2023.
- [8] A. Dias, J. Foleiss, and R. P. Lopes, "Reinforcement learning in tower defense," in *Conference on Videogame Sciences and Arts*, 2020.
- [9] J. T. Kristensen, A. Valdivia, and P. Burelli, "Statistical modelling of level difficulty in puzzle games," in *IEEE Conference on Games*, 2021.
- [10] S. Karimi, S. Asadi, F. Lorenzo, and A. H. Payberah, "Candyrl: A hybrid reinforcement learning model for gameplay," in *IEEE International Conference on Machine Learning and Applications*, 2022.
- [11] C. Politowski, Y.-G. Guéhéneuc, and F. Petrillo, "Towards automated video game testing: Still a long way to go," *arXiv preprint arXiv:2202.12777*, 2022.
- [12] L. Mugrai, F. Silva, C. Holmgård, and J. Togelius, "Automated playtesting of matching tile games," in *Conference on Games*. IEEE, 2019.
- [13] A. Sestini, J. Bergdahl, K. Tollmar, A. D. Bagdanov, and L. Gisslén, "Towards informed design and validation assistance in computer games using imitation learning," in *Conference on Games*, 2023.
- [14] J. Bergdahl, C. Gordillo, K. Tollmar, and L. Gisslén, "Augmenting automated game testing with deep reinforcement learning," in *Conference on Games*, 2020.
- [15] A. Sestini, L. Gisslén, J. Bergdahl, K. Tollmar, and A. D. Bagdanov, "Automated gameplay testing and validation with curiosity-conditioned proximal trajectories," *IEEE Transactions on Games*, 2022.
- [16] S. Abdelfattah, A. Brown, and P. Zhang, "Preference-conditioned pixel-based ai agent for game testing," in *Conference on Games*. IEEE, 2023.
- [17] Y. Shin, J. Kim, K. Jin, and Y. B. Kim, "Playtesting in match 3 game using strategic plays via reinforcement learning," *IEEE Access*, 2020.
- [18] Z.-J. Pang, R.-Z. Liu, Z.-Y. Meng, Y. Zhang, Y. Yu, and T. Lu, "On reinforcement learning for full-length game of starcraft," in *AAAI Conference on Artificial Intelligence*, 2019.
- [19] P. Massoudi and A. H. Fassihi, "Achieving dynamic ai difficulty by using reinforcement learning and fuzzy logic skill metering," in *International Games Innovation Conference*, 2013.
- [20] S. Liu, L. Chaoran, L. Yue, M. Heng, H. Xiao, S. Yiming, W. Licong, C. Ze, G. Xianghao, L. Hengtong *et al.*, "Automatic generation of tower defense levels using pcg," in *International Conference on the Foundations of Digital Games*, 2019.
- [21] J. Gillberg, J. Bergdahl, A. Sestini, A. Eakins, and L. Gisslén, "Technical challenges of deploying reinforcement learning agents for game testing in aaa games," in *Conference on Games*, 2023.
- [22] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.