

# Reward-Adaptive Iterative Discovery: A Case Study on Automated Game Testing for NHL26

**Anonymous authors**

Paper under double-blind review

## Abstract

1       Testing is a major effort for the gaming industry, requiring a significant part of development budget and people power. We present a case study on a development version of the ice hockey game *EA SPORTS NHL 26*, for which human playtesters test the goalie AI for behavioral exploits. To reduce the effort of re-testing the goalie AI after every game or behavior modification in the development phase, we propose Reward-Adaptive Iterative Discovery (RAID), a novel approach to automatically find exploits using an iterative Reinforcement Learning (RL) approach that trains a population of goal scoring agents. While previous approaches can already successfully find exploits, RL algorithms tend to overfit to a single solution. We introduce a simple extension on top of existing RL algorithms, such that they find multiple diverse high-quality solutions. For our first deployment of this approach, within a single experiment we were able to find six hockey scoring exploit strategies that were qualitatively similar to those that playtesters had found in hours-long manual testing sessions.

14 **Supplementary Video** This paper is accompanied by videos of the agents trained with our algorithm, RAID: <https://youtube.com/watch?v=59J11QFV9I0>

## 16 **1 Introduction**

17 With game worlds becoming vaster and game systems becoming more complex, the effort to test games grows likewise. To reduce repetitive aspects of game testing, prior work has investigated automating parts of it through the use of autonomous agents playing a game and reporting bugs (Ariyurek et al., 2019; De Woillemont et al., 2022; Bergdahl et al., 2020). While Reinforcement Learning (RL) agents can find individual exploits within a game (Baker et al., 2019), they tend to converge to one single “best” solution of solving a task. In this paper, we use *EA SPORTS NHL 26* (short NHL) as our test-case. In particular, we aim to test the goalie AI of the game, training a forward agent to find high chance scoring strategies, which could represent potential exploits. Our results indicate that standard RL algorithms tend to collapse to a small set of high-reward behaviors, rather than exploring a larger, more diverse set. Thus, exploit discovery becomes sequential: developers must fix one issue before retraining the agent to find others. Since this slows down the exploit finding process, we aim to develop an algorithm that can find multiple potential exploits without human intervention.

30 Prior work on RL for game testing has proposed methods to increase coverage of a game environment (Gordillo et al., 2021; Sestini et al., 2022). However, the added bias towards exploration comes at the cost of a loss in optimality with respect to the original goal of the agent. For our use case, we are explicitly looking for high-performing scoring strategies and we therefore require a method that achieves high performance with respect to the scoring chance, while also finding multiple diverse solutions without human intervention.

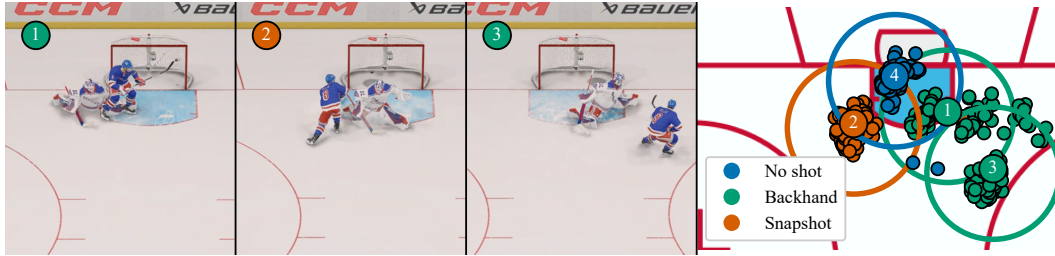


Figure 1: **Left: Last frame before the goal of first 3 scoring strategies learned by agents trained with RAID.** The controlled forward player wears blue, the goalie white. The strategies are learned in an iterative fashion, with a reward function enforcing each new strategy’s shot position to be at least 2 meters away from all previous strategies using the same shot type. Numbers indicate the algorithm iteration. Colors indicate the shot type. **Right: Shot locations and types of strategies found by first 4 iterations of RAID.** Small dots represent shot location and type of 100 goals after convergence of each iteration. Large circles represent 2 meter radius around average position of those goals inside which no reward is given for all succeeding iterations of RAID if scoring with the same shot type. Shots by the agent of iteration 3 are for example outside of the 2 meter radius of iteration 1, since the two agents use the same shot type.

36 The field of quality diversity investigates how to find multiple diverse solution that are high-  
 37 performing (Zahavy et al., 2023). While they have reported impressive results, existing solutions  
 38 are often brittle in complex domains, one reason being a moving diversity target learned in parallel  
 39 to the policies (Leon et al., 2024). We aim for playtesters to independently use our tool and therefore  
 40 want the tool to be as simple as possible to require little RL expertise.

41 Following these requirements, we propose Reward-Adaptive Iterative Discovery (RAID), an ap-  
 42 proach to find a set of high-quality, diverse solutions by iteratively training a population of agents.  
 43 Figure 1 shows a brief overview of the results of the approach. We enforce diversity through not  
 44 rewarding agents for strategies that are similar to those found by any previously trained agent. By  
 45 measuring diversity over sequentially trained agents, leaving the diversity measure fixed during an  
 46 individual agent’s training time, we increase robustness compared to previous work. We define di-  
 47 versity in an intuitive, domain-specific way, such that the diversity measure can be interpreted and  
 48 adjusted by non-RL practitioners, facilitating the use by playtesters to extend their toolset.

## 49 2 Related Work

50 We present related work in two research directions connected to our approach: (1) automated  
 51 playtesting, specifically approaches based on reinforcement learning, and (2) the field of behavior  
 52 diversity, which motivates our approach.

### 53 2.1 Automated Playtesting

54 Contrary to recent work that apply reinforcement learning to train agents to master a game and  
 55 surpass human players (Vinyals et al., 2019; Wurman et al., 2022), research on automated playtesting  
 56 is generally less concerned with maximizing the same success metrics as human players, and more  
 57 focused on having agents play a game in various ways that reveal insights on potential issues in the  
 58 game. One direction for this is to maximize the coverage of states an agent visits in a game. A  
 59 notable example is the work by Gordillo et al. (2021), that uses a count-based exploration reward  
 60 scheme to promote visitation of less seen states in a game map. Another direction aims to train agents  
 61 that play games with behavior similar to human players, for example to visualize the interactions a  
 62 game designer can expect from users or to balance the effectiveness of different strategies within a  
 63 game (Holmgård et al., 2019; De Woillemont et al., 2022; Devlin et al., 2021). Sestini et al. (2022)

64 combine both directions by exploring strategies similar to those of human players with an additional  
65 intrinsic reward for exploring new states.

66 While prior work achieves good results in terms of game state coverage and player-like behavior,  
67 we aim to investigate exploits that are closely tied to high performance with respect to the original  
68 human success metric of a game, in our case scoring goals with a high percentage in a hockey game.  
69 We therefore aim to develop an approach to increase the diversity in behaviors of the trained agents  
70 while preserving the performance with respect to the game’s success metric as much as possible.

## 71 2.2 Diversity

72 The field of quality diversity aims to generate diverse populations of high-performing solutions  
73 through evolutionary algorithms (Lehman & Stanley, 2011; Mouret & Clune, 2015; Cully et al.,  
74 2015). More recent work brings the concept of diversity into the field of reinforcement learning  
75 through intrinsic rewards. The work by Gregor et al. (2017), for example, introduces an intrinsic  
76 reward based on optimizing the mutual information between options – i.e. closed-loop policies for  
77 taking action over a period of time – and the options’ final states, using entropy optimization to  
78 maximize diversity across options and distinction of individual options. Similarly Eysenbach et al.  
79 (2019) optimize diversity over all visited states and show that this form of unsupervised discov-  
80 ery can serve as an effective pretraining mechanism for reinforcement learning. DOMiNO defines  
81 diversity in the space of state-action occupancy (Zahavy et al., 2023), represented by successor  
82 features (Barreto et al., 2017). It combines the intrinsic diversity reward and an extrinsic reward  
83 specific to the environment with Lagrange multipliers to control the trade off between diversity and  
84 optimality with respect to the extrinsic reward. While DOMiNO can train multiple diverse agents  
85 in parallel this way, Leon et al. (2024) highlight that the approach fails to learn diverse behaviors in  
86 complex domains and can suffer from instability due to the approach learning the successor feature  
87 representations in parallel to the agents’ policies.

88 Due to this instability and our aim for a stable algorithm that can be used by non-RL practitioners,  
89 we train a group of agents in a simple, sequential fashion, by masking strategies similar to those  
90 found by a prior agent from the reward function of all subsequent agents. The diversity measure-  
91 ment introduced by this method of reward masking is static for a specific agent within a single  
92 training iteration and not any more complex in practice than simple static reward shaping in stan-  
93 dard reinforcement learning (Ng et al., 1999). Similar to prior work by Gregor et al. (2017), we  
94 define diversity in terms of the final state of an agent. Rather than maximizing the distance between  
95 states, we however introduce diversity via hard constraints on state difference, enforced through  
96 reward masking of prior behaviors. This allows users of RAID to easily update the algorithm to  
97 their desired level of diversity, for example in the later introduced diversity criterion for NHL by  
98 increasing the area used for the shot similarity criterion to a 4 m radius. In comparison, the diversity  
99 hyperparameters of prior work, for example relative to the maximum achievable reward in the work  
100 by Zahavy et al. (2023), are significantly harder to interpret, especially for non-RL practitioners.

## 101 3 Methodology

102 Our goal is to create an RL-based approach to autonomously find multiple diverse, high chance  
103 scoring strategies for the NHL environment without any human intervention in between iterations.  
104 We first outline the base RL setup to find a single exploit within the NHL environment. We then  
105 describe RAID, a novel approach to satisfy the posed requirements. For the foundations of reinforce-  
106 ment learning, such as the concepts of reward, policy, and Q-value function, we refer the reader to  
107 the work by Sutton & Barto (2018).

### 108 3.1 Base RL Setup

109 **Action Space.** We aim to train a hockey forward player agent that learns high chance scoring  
 110 strategies in a one-versus-goalie setting, similar to those that a human player could find by exploiting  
 111 weaknesses in the goalie AI system. To make sure that the found strategies are executable by a  
 112 human, we provide the agent with a selection of the same actions available to human players on a  
 113 game controller:

- 114 • Discrete actions: the 4 face buttons and the left trigger button that execute so-called *dekes*, i.e.,  
 115 decoy movements to draw the goalie out of position.
- 116 • Continuous actions: the two sticks for player skating and hockey stick control, each represented  
 117 by continuous  $x$  and  $y$  axes.

118 Through the use of those actions, the agent can perform different kinds of shot types – backhand,  
 119 wrist shot, etc. – which we will later consider for defining shot diversity. Furthermore, we only let  
 120 the agent act every 5 frames and smooth the stick actions using an exponential moving average with  
 121 a smoothing factor of 0.2, to make sure that the agent does not leverage any super-human movement  
 122 capability or reaction time.

123 **Observation Space.** We provide the agent with the normalized position, velocity, and orientation  
 124 of puck, net, and goalie relative to the agent, as well as a stack of the 8 last performed actions.

125 **Reward Function.** We reward the agent for every goal. To simplify learning, we additionally  
 126 include shaping rewards for the puck getting closer to the goal, as well as for the scoring chance of  
 127 every shot taken, based on an NHL-internal calculation.

128 **Algorithm & Architecture.** While from an interface perspective all RL algorithms that use a  
 129 reward function are technically compatible with RAID, RAID’s effectiveness depends on the base  
 130 algorithm’s ability to explore various strategies – even when previously found strategies are no  
 131 longer rewarded – before converging to the best one. We therefore use Soft Actor-Critic ([Haarnoja](#)  
 132 [et al., 2018](#)), since it features a robust exploration mechanism through entropy maximization, ex-  
 133 ploring more in less known states. Furthermore, Soft Actor-Critic has a high sample efficiency,  
 134 which is crucial for environments with slow data generation, such as full-scale games. To support  
 135 both discrete and continuous actions, we use the architecture described by [Delalleau et al. \(2019\)](#)  
 136 for 2 Q-value functions, as well as a mixed-action policy, each with 5 hidden layers with 512 units,  
 137 allowing for a maximum of one discrete action at a time. When we use the term *agent*, we refer to  
 138 both the policy and Q-value functions used for training and inference. For efficient training, we stop  
 139 training once an agent either reaches a 90% scoring chance or does not improve its performance for  
 140 50k training steps.

### 141 3.2 Reward-Adaptive Iterative Discovery

142 We outline the general, domain-agnostic RAID method as a pseudo algorithm in Appendix A. The  
 143 algorithm extends a standard reinforcement learning setup by training multiple agents sequentially,  
 144 aiming for each agent to come up with a new strategy for maximizing the reward of the environment.  
 145 For NHL we define a strategy  $z$  as the shot type and shot position at the end of an episode, the shot  
 146 position meaning the last position of the puck before the puck is shot or, if the puck is not shot, the  
 147 last puck position before a goal. We start by training a first agent with a standard reinforcement  
 148 learning setup. After training of the first agent converges, we evaluate the agent to capture the  
 149 strategy of its highest performing checkpoint. For the NHL game, we evaluate the agent until it  
 150 scores 100 goals and describe the strategy  $z$  as the average shot position and the most common shot  
 151 type over those 100 goals. We store this strategy in a list of previous strategies  $Z_{\text{prev}}$ . We then  
 152 continue sequentially, training each next agent the same way, but for each new iteration modify the  
 153 reward function to not reward strategies similar to those found by any previous agent. Based on  
 154 feedback from the NHL development team, we define  $\text{similar}(z, Z_{\text{prev}})$  as the shot position of a

155 strategy being within a 2 m radius of a previous strategy that used the same shot type, i.e. a strategy  
 156  $z$  is considered novel or sufficiently diverse with respect to a set of existing strategies  $Z_{\text{prev}}$  if:

$$157 \forall z' \in Z_{\text{prev}} \quad z_{\text{shottype}} \neq z'_{\text{shottype}} \vee l^2(z_{\text{shotpos}}, z'_{\text{shotpos}}) > 2 \text{ m.}$$

158 For NHL, we skip adding strategies with a scoring chance of less than 10% after convergence, since  
 159 we are interested in high-performing strategies and convergence at this low performance potentially  
 160 indicates that the algorithm got stuck in a local optimum, meaning that future iterations starting with  
 161 a different random seed could potentially find more performant solutions for the same shot position  
 162 and type. For NHL, we use the details from Section 3.1 for the base reward function, training  
 163 algorithm, and convergence criterion.

164 While the definition of strategy similarity is highly domain specific, we argue that this is exactly  
 165 what makes the method suitable for game development. While previous work defines diversity in  
 166 a more general way that can be used across domains, for example via successor features (Zahavy  
 167 et al., 2023), we argue this makes it significantly harder for a playtester without RL knowledge  
 168 to tune diversity to a level fitting their domain, given that the work by Zahavy et al. (2023), for  
 169 example, requires the diversity level to be expressed with respect to the maximum achievable reward.  
 170 Furthermore, having an explicit list of previous strategies allows playtesters to warm start training  
 171 using a list of manually defined strategies, instead of an empty set, for example using high scoring  
 172 strategies that are known and desired by the game, or exploits that are known but the game studio  
 173 could not fix yet.

## 174 4 Experiments

175 We compare RAID to a naive baseline to demonstrate the diversity in strategies found by the ap-  
 176 proach. This is most notably supported by RAID finding 6 exploit strategies that were previously  
 177 also found by human playtesters, but had not been fixed on the outdated version of the game that we  
 178 conducted the experiments on. We furthermore test the limitations of the algorithm with respect to  
 179 the number of high chance scoring strategies it finds before failing to find any new strategies.

### 180 4.1 Experimental Setup

181 We run all experiments on a single PC with an *NVIDIA GeForce RTX 4090* GPU and an *AMD*  
 182 *Ryzen Threadripper PRO 7975WX* CPU. We use a pre-release development version of *EA SPORTS*  
 183 *NHL 26*, additionally reducing rendering to a minimum to speed up the simulation. The simulation  
 184 includes a single goalie, controlled by the game’s AI we want to test, as well as a forward player,  
 185 controlled by RL. For both training and evaluation we slightly vary the initial position of the forward  
 186 around the centerline to make the approach more robust to variations in the setup and to help with  
 187 exploration.

### 188 4.2 Naive Baseline

189 As a baseline, we run the base RL setup from Section 3.1 for 20 independent iterations, each with  
 190 a different random seed, changing the initialization of the policy and value networks, as well as  
 191 the random states of the game environment. Each iteration uses the full reward function, assigning  
 192 rewards for all goals, independent of shot position or type, without masking previous strategies.  
 193 The experiments take between 16 min and 152 min to converge. RL algorithms aim to find the  
 194 best solution to a problem with respect to a reward function. In environments with a set of clearly  
 195 superior strategies, returning a higher reward, an optimal algorithm therefore always converges to  
 196 this same set of solutions, if not otherwise incentivized. The left of Figure 2 shows how all of the  
 197 20 independent iterations end up with either one of two patterns, a snapshot from the left side of the  
 198 goal or a backhand shot from the right side.

199 Assuming that a playtester classifies the found strategies as unwanted exploits, they have to improve  
 200 the goalie AI behavior to mitigate those exploits before the baseline approach is likely to find any

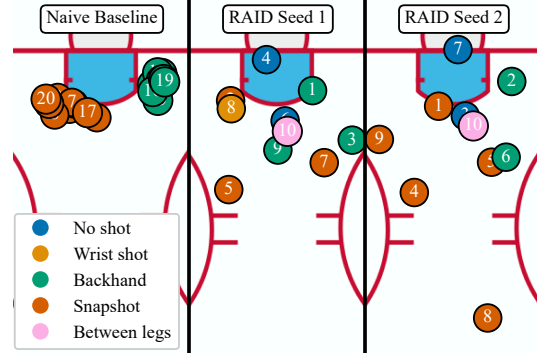


Figure 2: The average shot position and most common shot type after convergence for 20 independent repetitions of the baseline experiment, as well as two experiments using RAID started with different random seeds ran for 10 iterations each. The color represents the most common shot type of an agent and the number represents the iteration of the algorithm that leads to the shot strategy. The naive baseline shows that a standard reinforcement learning algorithm keeps finding the same small set of solutions when repeated multiple times, with the algorithm finding 9 snapshot strategies and 11 backhand strategies, all from similar positions. The two experiments using RAID not only find the two strategies found by the baseline, but additionally also find 8 more strategies that adhere to our definition of diversity, maintaining a distance of at least 2 m between shot strategies that use the same shot type.

201 additional ones. The naive RL approach therefore adds considerable overhead to the exploit finding  
 202 process by requiring human intervention between iterations. If on the contrary the playtester con-  
 203 siders the found strategies to be a desired way of scoring in the game, they would need to find a  
 204 way for the reward function to distinguish between exploits and desired scoring strategies, such that  
 205 the agent ignores the desired strategies and finds novel exploits. Exploits are highly unique in their  
 206 appearance and therefore hard to define before first seeing them.

### 207 4.3 Reward-Adaptive Iterative Discovery

208 We run RAID starting with two different initial random seeds, influencing the random initialization  
 209 of the agents of all iterations, for  $N = 30$  and  $N = 10$  iterations respectively. The 30 iteration  
 210 experiment takes 48 h and the 10 iteration experiment takes 14 h for all iterations to converge. Fig-  
 211 ure 2 shows how the strategies found by RAID adhere to the diversity criteria we set for NHL: while  
 212 the average shot positions of strategies with the same shot type are all at least 2 m away from each  
 213 other, for example the three backhand shot strategies in green for seed 1, the algorithm still finds  
 214 strategies with close shot position but different shot types, e.g. strategies 8 and 2 for seed 1, sharing  
 215 almost the same position. This way the algorithm finds a set of diverse strategies without any human  
 216 intervention in between iterations to fix prior exploits, as would be required for the naive baseline.

217 As an illustrative example of the usefulness of RAID for playtesting, when we ran RAID for the  
 218 first time and presented the results to the development team, 6 of the found strategies matched with  
 219 exploits previously found by human playtesters. Since we did however use an outdated version of  
 220 the game for our experiments, those exploits were still available for the agents to find. This shows  
 221 that the agents are able to find exploits qualitatively similar to those found by human playtesters,  
 222 while only requiring human supervision after the algorithm has found a set of candidates for a  
 223 human expert to review. To not reveal any information that allows players to further exploit game  
 224 mechanisms, we do not reveal which of the found strategies match with those deemed exploits.

225 **Pre-shot Diversity.** Figure 3 shows 10 trajectories each for a selection of iterations of seed 1 of  
 226 RAID. While we only define diversity in terms of shot position and type, the visualized trajectories  
 227 demonstrate how this diversity definition implicitly generates diversity in states before the shot as  
 228 well. This shows that the agents take into account player direction, as well as orientation and stick

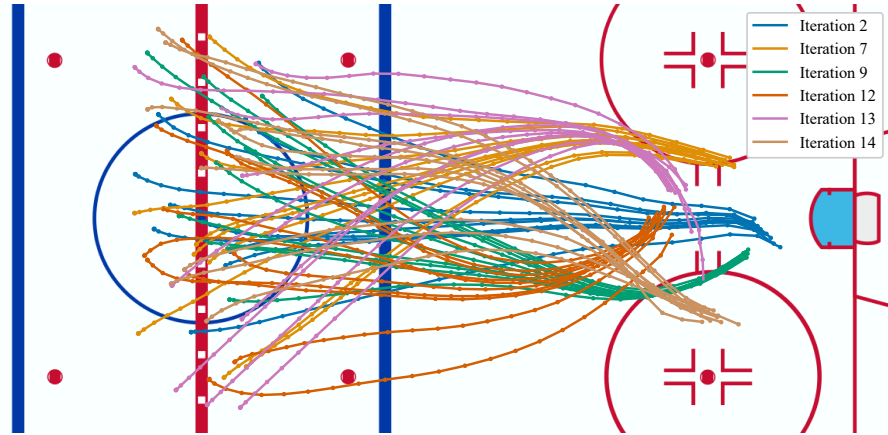


Figure 3: Movement of the agents of a hand-picked selection of iterations of seed 1 of RAID for 10 different random starting positions each. This demonstrates that while diversity is only defined in terms of shot type and position of the shot, the positions before the shot, are also implicitly impacted to create a high scoring chance. For example, while iterations 12 and 13 have almost the same shot position, they approach the position from the opposite side to facilitate different shot types.

229 position, which are visible in the supplementary videos, when learning scoring strategies. While  
 230 trajectories found by iterations 12 and 13, for example, have similar shot positions, they require  
 231 significantly different skating trajectories to score with their respective shot types, backhand and  
 232 snapshot.

233 **No-shot Goals.** To our surprise, the agent also finds strategies to score without shooting, marked  
 234 in dark blue in Figure 2 and visible in the supplementary video. This initially led to a bug in our  
 235 setup, since previous shot strategies were only recorded when the agent actually shot the puck. This  
 236 further underlines the difficulty of defining exploits before first seeing them, outlined in Section 4.2.

237 **Re-testing after Behavior Update.** While most of the strategies found in seed 1 and seed 2  
 238 resemble each other, they appear in different order and do not match completely – for instance,  
 239 strategies found in iteration 8 of seed 1 and iteration 8 of seed 2 are not found by the other seed.  
 240 RL algorithms have a high variance in their exploration process and therefore converge to different  
 241 solutions given the existence of multiple similarly performant solutions, as we also demonstrate  
 242 with the two dominant solutions the naive baseline finds. This however means that if a playtester  
 243 attempts to fix an exploit found by a first run of RAID, the same exploit not appearing in a second  
 244 run of RAID after the fix is not a sufficient guarantee that the fix was successful. Similarly, the  
 245 policies trained in the first run can not be used for this purpose, since RL is highly overfitting to the  
 246 exact dynamics of the environment and can therefore be easily thrown off by slight changes in the  
 247 goalie’s behavior, without proving that the exploit is fixed. Playtesters therefore still have to validate  
 248 that an exploit is indeed fixed by imitating the strategy found by the approach before the fix.

249 **Iteration Limit.** Figure 4 shows the training progress of the 30 iterations of seed 1. The first  
 250 10 iterations all reach performances above the 10% goal rate cut-off criterion before triggering any  
 251 of the early stopping criteria. The longer the experiment goes on, the more iterations do not make  
 252 the cut, with iterations 11, 15, 17, 20, 21, 23 – 25 and 27 – 29 all converging at goal rates below  
 253 the threshold and therefore not being added to the previous strategies list ( $Z_{prev}$ ). We manually stop  
 254 the search after 30 iterations, since we only find 3 valid strategies in the last 11 iterations, reducing  
 255 the likelihood of discovering further interesting strategies in practical time. More generally, limit  
 256 allows for a trade-off between run time and the expected number of diverse solutions found, with  
 257 the optimal choice being highly domain-specific and dependent on the diversity criterion. In NHL  
 258 for example, if we did not differentiate the shot type and only used the position as discriminative  
 259 feature, we would expect to find less solutions over all, reducing the optimal search time.

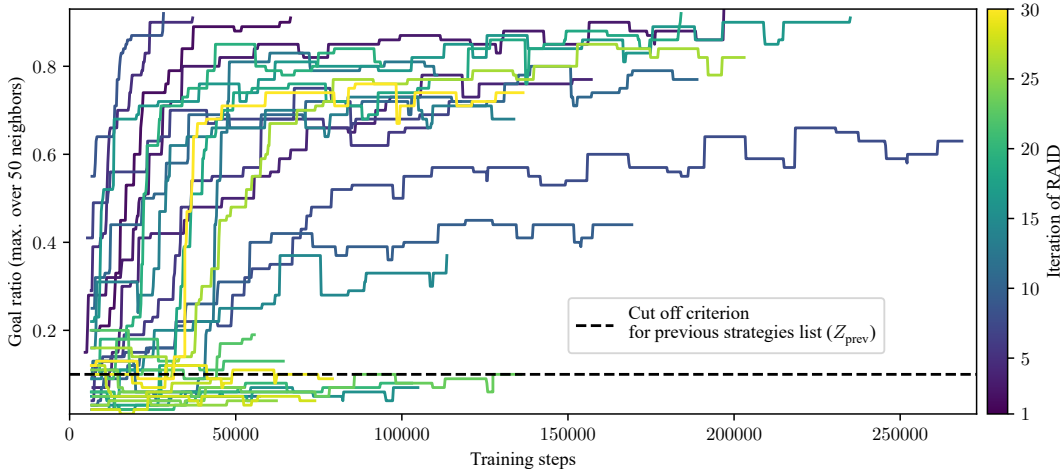


Figure 4: Training progress of seed 1 of RAID running for 30 iterations, represented by the goal rate of the agents during training episodes. For this figure, we smooth the goal rate using a maximum over 50 neighbors to maintain information on whether an iteration reaches the goal rate threshold of 10% at any point in training. The coloring shows that later iterations of RAID, which are represented in brighter colors, are less likely to discover high chance scoring strategies, because most of those strategies were already found by previous iterations, represented in darker colors.

260 **Multimodality.** Our approach assumes that a 2 m radius around the mean shot position sufficiently  
 261 covers the shots by a single agent. We did however find that especially in later training iterations  
 262 this is not the case. In Appendix B we highlight the issue in more detail.

## 263 5 Conclusion, Limitations, and Future Work

264 We introduce Reward-Adaptive Iterative Discovery, a method to find multiple diverse high chance  
 265 scoring strategies without human intervention in between iterations. The algorithm extends a stan-  
 266 dard reinforcement learning setup by training multiple agents sequentially, aiming for each agent  
 267 to come up with a new strategy by not rewarding agents for strategies similar to those of previous  
 268 iterations. As a domain expert’s review of the learned strategies shows, RAID finds exploits similar  
 269 to those identified by human playtesters. We furthermore argue that defining diversity with respect  
 270 to NHL allows game testers to better follow and adjust the method, compared to prior work that  
 271 requires RL expertise to interpret the hyperparameters for tuning diversity.

272 **Limitations.** While our method finds potential exploits autonomously, it still requires human  
 273 supervision post training, since telling high chance scoring strategies intended by game designers  
 274 and actual exploits apart is not covered by the method. Furthermore, due to the high variance in RL’s  
 275 exploration process, our method is not sufficient to determine whether a found exploit has been fixed  
 276 successfully by an engineer, relying on playtesters for this verification task. Our method’s runtime  
 277 grows (at least) linearly with the number of strategies we aim to find. While previous work trains  
 278 multiple diverse behaviors in parallel (Zahavy et al., 2023), we deliberately forego this option to  
 279 keep our method simple and therefore facilitate its use by playtesters not familiar with RL.

280 **Future Work.** We design parts of the RAID method specifically for the use on NHL, though the  
 281 base idea outlined in Algorithm 1 could be applicable to other domains as well. While diversity is  
 282 somewhat intuitive to define for the domain of scoring goals in NHL, future work could research how  
 283 diversity can be defined in more complex domains, including multiplayer scenarios or even agents  
 284 outputting text as actions. While previous work that defines diversity in a domain-agnostic way often  
 285 requires RL knowledge to interpret and adjust hyperparameters, it would be interesting to research  
 286 definitions of diversity that are domain-agnostic but still interpretable by non-RL practitioners, for  
 287 example by defining the minimum distance of trajectories used in RAID in a more general way.

288 **References**

- 289 Sinan Ariyurek, Aysu Betin-Can, and Elif Surer. Automated video game testing using synthetic and  
290 humanlike agents. *IEEE Transactions on Games*, 13(1):50–67, 2019.
- 291 Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew, and Igor  
292 Mordatch. Emergent tool use from multi-agent autocurricula. In *International conference on*  
293 *learning representations*, 2019.
- 294 André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P Van Hasselt,  
295 and David Silver. Successor features for transfer in reinforcement learning. *Advances in neural*  
296 *information processing systems*, 30, 2017.
- 297 Joakim Bergdahl, Camilo Gordillo, Konrad Tollmar, and Linus Gisslén. Augmenting automated  
298 game testing with deep reinforcement learning. In *2020 IEEE Conference on Games (CoG)*, pp.  
299 600–603. IEEE, 2020.
- 300 Antoine Cully, Jeff Clune, Danesh Tarapore, and Jean-Baptiste Mouret. Robots that can adapt like  
301 animals. *Nature*, 521(7553):503–507, 2015.
- 302 Pierre Le Pelletier De Woillemont, Rémi Labory, and Vincent Corruble. Automated play-testing  
303 through rl based human-like play-styles generation. In *Proceedings of the AAAI Conference on*  
304 *Artificial Intelligence and Interactive Digital Entertainment*, volume 18, pp. 146–154, 2022.
- 305 Olivier Delalleau, Maxim Peter, Eloi Alonso, and Adrien Logut. Discrete and continuous action  
306 representation for practical rl in video games. *arXiv preprint arXiv:1912.11077*, 2019.
- 307 Sam Devlin, Raluca Georgescu, Ida Momennejad, Jaroslaw Rzepecki, Evelyn Zuniga, Gavin  
308 Costello, Guy Leroy, Ali Shaw, and Katja Hofmann. Navigation turing test (ntt): Learning to  
309 evaluate human-like navigation. In *International Conference on Machine Learning*, pp. 2644–  
310 2653. PMLR, 2021.
- 311 Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need:  
312 Learning skills without a reward function. In *International Conference on Learning Representa-*  
313 *tions*, 2019. URL <https://openreview.net/forum?id=SJx63jRqFm>.
- 314 Camilo Gordillo, Joakim Bergdahl, Konrad Tollmar, and Linus Gisslén. Improving playtesting  
315 coverage via curiosity driven reinforcement learning agents. In *2021 IEEE Conference on Games*  
316 *(CoG)*, pp. 1–8. IEEE, 2021.
- 317 Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *Interna-*  
318 *tional Conference on Learning Representations, Workshop Track*, 2017.
- 319 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
320 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*  
321 *ence on machine learning*, pp. 1861–1870. Pmlr, 2018.
- 322 Christoffer Holmgård, Michael Cerny Green, Antonios Liapis, and Julian Togelius. Automated  
323 playtesting with procedural personas through mcts with evolved heuristics. *IEEE Transactions on*  
324 *Games*, 11(4):352–362, 2019. DOI: 10.1109/TG.2018.2808198.
- 325 Joel Lehman and Kenneth O. Stanley. Evolving a diversity of virtual creatures through novelty  
326 search and local competition. In *Proceedings of the 13th Annual Conference on Genetic and*  
327 *Evolutionary Computation*, GECCO ’11, pp. 211–218, 2011. ISBN 9781450305570. DOI: 10.  
328 1145/2001576.2001606. URL <https://doi.org/10.1145/2001576.2001606>.
- 329 Borja G Leon, Francesco Riccio, Kaushik Subramanian, Peter R Wurman, and Peter Stone. Dis-  
330 covering creative behaviors through duplex: Diverse universal features for policy exploration.  
331 *Advances in Neural Information Processing Systems*, 37:49625–49648, 2024.

- 332 James B McQueen. Some methods of classification and analysis of multivariate observations. In  
 333 *Proc. of 5th Berkeley Symposium on Math. Stat. and Prob.*, pp. 281–297, 1967.
- 334 Jean-Baptiste Mouret and Jeff Clune. Illuminating search spaces by mapping elites. *arXiv preprint*  
 335 *arXiv:1504.04909*, 2015.
- 336 Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations:  
 337 Theory and application to reward shaping. In *Icml*, volume 99, pp. 278–287. Citeseer, 1999.
- 338 Alessandro Sestini, Linus Gisslén, Joakim Bergdahl, Konrad Tollmar, and Andrew David Bag-  
 339 danov. Automated gameplay testing and validation with curiosity-conditioned proximal trajec-  
 340 tories. *IEEE Transactions on Games*, 16(1):113–126, 2022.
- 341 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press,  
 342 second edition, 2018. URL [http://incompleteideas.net/book/the-book-2nd.](http://incompleteideas.net/book/the-book-2nd.html)  
 343 [html](http://incompleteideas.net/book/the-book-2nd.html).
- 344 Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Juny-  
 345 oung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster  
 346 level in starcraft ii using multi-agent reinforcement learning. *nature*, 575(7782):350–354, 2019.
- 347 Peter R Wurman, Samuel Barrett, Kenta Kawamoto, James MacGlashan, Kaushik Subramanian,  
 348 Thomas J Walsh, Roberto Capobianco, Alisa Devlic, Franziska Eckert, Florian Fuchs, et al. Out-  
 349 racing champion gran turismo drivers with deep reinforcement learning. *Nature*, 602(7896):223–  
 350 228, 2022.
- 351 Tom Zahavy, Yannick Schroecker, Feryal Behbahani, Kate Baumli, Sebastian Flennerhag, Shaobo  
 352 Hou, and Satinder Singh. Discovering policies with DOMiNO: Diversity optimization main-  
 353 taining near optimality. In *The Eleventh International Conference on Learning Representations*,  
 354 2023. URL <https://openreview.net/forum?id=kjkdzBW3b8p>.

## 355 A Pseudo algorithm

---

### Algorithm 1 Reward-Adaptive Iterative Discovery

---

**Input:**  $N, Z_{\text{prev}}$  ▷  $Z_{\text{prev}}$ : previous strategies, can be  $\emptyset$  or warm-started  
**Output:**  $Z_{\text{prev}}$   
**for**  $i = 1$  to  $N$  **do**  
    $r_i(z) \leftarrow r_{\text{base}}(z)$  **IF** !similar( $z, Z_{\text{prev}}$ ) **ELSE** 0 ▷  $z$ : strategy,  $r(z)$ : reward function  
   Initialize agent  $\pi_i$   
   **while** NOT converged( $\pi_i$ ) **do**  
      train( $\pi_i, r_i$ )  
   **end while**  
    $z_i \leftarrow \text{evaluate}(\pi_i)$   
    $Z_{\text{prev}} = Z_{\text{prev}} \cup \{z_i\}$   
**end for**

---

## 356 B Multimodality of shot distribution

- 357 Figure 5 shows the shot position of 100 goals after training converged for the 8th iteration of seed 2.  
 358 Around one third of all shot positions lie outside of the 2 m radius around the average of all shot  
 359 positions, meaning that in succeeding training iterations of RAID the exact same shots will still be  
 360 rewarded. This highlights a potentially bigger problem around the introduced diversity criterion for  
 361 NHL: in a more extreme hypothetical case, a single agent can learn to shoot from two positions that

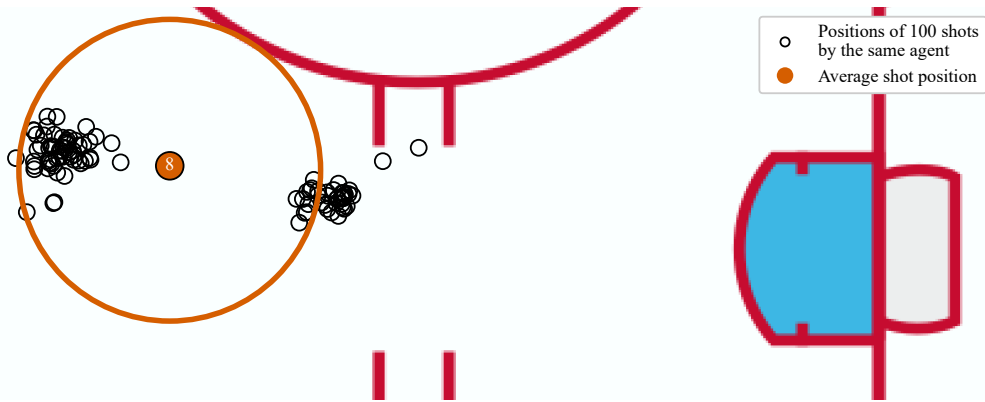


Figure 5: 100 shot positions of the agent in iteration 8 of seed 2 of RAID, as well as the average of those shot positions. In succeeding iterations, shots within 2 m of that average position are no longer rewarded if they use the same shot type as the most common shot type among those 100 shots. The agent learns a bi-modal behavior, shooting from either one of two shot positions depending on its spawn position. This leads to around 1/3 of shots not being within the radius. Succeeding iterations of the algorithm are therefore still rewarded for shots similar to this part of the behavior, meaning that the exclusion mechanism based on the radius around the average does not completely exclude this scoring strategy.

362 are more than 4 meters apart, going for either of the two shot positions at a rate of 50% depending  
 363 on where the agent is spawned. This way, most of the shots of the agent could be outside of the  
 364 2 m radius around the average of the shots' positions, meaning that future iterations could still learn  
 365 the exact same behavior without it being excluded. This motivates further investigations into more  
 366 elaborate representations of a policy's strategy, for example through clustering methods such as  
 367 k-means (McQueen, 1967).