

1 We thank all of the reviewers for their valuable feedback. We appreciate that **Reviewer 4** finds the theoretical connection  
2 to decision theory both “sound” and “clear and valuable” and that **Reviewer 3** states that “the formulations are correct.”  
3 However, we acknowledge that the clarity and organization of the theory could be improved. To address this, we will  
4 explain the major theoretical results more clearly in the main body of the text, and include a sketch of the proof of  
5 Theorem 2 as suggested by **Reviewer 4**. We will incorporating **Reviewer 2**’s suggestions for improving the clarity of  
6 the goals and assumptions of the theoretical results. We will also provide additional experiments in the appendix, which  
7 test new methods for computing regret (for **Reviewer 1**) and assess whether a curriculum can be generated without  
8 training the protagonist using regret (for **Reviewer 4**).

9 **Reviewer 2** points out that Theorem 2 shows that PAIRED recovers the minimax regret policy when the adversary and  
10 the antagonist successfully coordinate, but it is unclear this condition holds. However, coordination is not required for  
11 PAIRED to work. If the policies reach a Nash Equilibrium, the protagonist performs at least as well as the antagonist in  
12 every parameterization, since the protagonist could guarantee 0 payoff by using the antagonist’s strategy and otherwise  
13 the adversary would choose a distribution on which the protagonist performs relatively worse. With a capable antagonist  
14 the protagonist would learn the minimax regret policy, even without the adversary and the antagonist coordinating. We  
15 will add this explanation to the paper.

16 **Reviewer 2** asks for clarification of the conditions of Theorem 1: there are two outcomes, Success and Failure, and the  
17 range of rewards given for Success does not overlap with the range of rewards given for Failure. Specifically, Successful  
18 outcomes give rewards in some range  $[S_{min}, S_{max}]$ , and Failure outcomes give rewards some in range  $[F_{min}, F_{max}]$   
19 such that  $F_{min} \leq F_{max} < S_{min} \leq S_{max}$ , and  $S_{max} - S_{min} < S_{min} - F_{max}$  and  $F_{max} - F_{min} < S_{min} - F_{max}$ . We  
20 thank **Reviewer 2** for bringing this ambiguity to our attention and will update the paper accordingly.

21 **Reviewer 2** also offered the valuable advice to improve clarity by offering a “specific example/ problem that would  
22 be easily solved by defining it as a UED”. One example is training a robot in simulation to pick up objects from a  
23 bin in a real-world warehouse. There are many possible configurations of objects, including objects being stacked  
24 on top of each other. We may not know *a priori* the typical arrangement of the objects, but can naturally describe  
25 them as parameterizations of a simulated environment. We can provide a curriculum of training configurations with  
26 either domain randomization, minimax, or PAIRED. We will introduce this example in the introduction and reference it  
27 throughout the formal introduction of UED to improve the clarity of the formalism.

28 **Reviewer 2** points out that “clarity for UED could be improved” with “a clearer description saying we have X and  
29 want to achieve Y”. In our setting we are given a class of training environments, and our goal is to construct a policy  
30 which performs well across a large set of these environments. We evaluate with a set of specific environments used as  
31 transfer tasks. To train such a policy, we start with an initially random policy, generate environments which are tuned to  
32 help it learn, train the policy on those environments, and repeat until convergence or until we have exceeded available  
33 computational resources. We call the problem of choosing how to generate these environments UED. The choice of  
34 how to solve the UED problem affects both the curriculum it generates and how the policy prioritizes performance in  
35 different environments at convergence. We will add phrasing to this effect in the introduction of the paper.

36 **Reviewer 1** points out that although using the ‘relative regret’ between the protagonist and an imperfect antagonist leads  
37 to an effective curriculum, it may be an inaccurate estimate of the worse-case regret during training. They add, “There  
38 are potentially many choices for how to differentially train the protagonist/antagonist policies”. We are conducting  
39 additional experiments empirically investigating alternative methods for estimating regret: 1) we use a population  
40 of agents and compute regret as the difference between the highest performing agent and the population average, 2)  
41 we relabel the the current best-performing agent as the antagonist. We will include the results of these additional  
42 experiments in the appendix.

43 **Reviewer 4** asked whether it is possible to decouple generating a curriculum of environments using regret, and having  
44 the agent itself optimize regret. We have experimented with training the protagonist with environmental reward and not  
45 regret and find that it still performs well; we will provide these results in the appendix. In addition, the performance  
46 of PAIRED in the random 50 blocks environment indicates that it would perform well as a policy for the domain  
47 randomization setting, even though it was designed for a different objective.

48 **Reviewer 3** - we will provide further discussion on scenarios where the method is not applicable, including settings in  
49 which you already have an accurate model of the test environments you expect to encounter.

50 We will improve the Broader Impact statement according to the suggestions of **Reviewer 2**, making it more specific and  
51 focusing on PAIRED’s ability to train more robust agents that “cover corner cases that may happen in the real world”.  
52 We will correct all typos identified by the reviewers, including the caption of Figure 3 which mistakenly refers to a Lava  
53 environment. We appreciate your detailed feedback.