

1 We thank the reviewers for their feedback. In particular, (R1, R3, R5) asked many specific and thoughtful questions,
2 with (R1) marking questions by priority. Thank you.

3 We’re glad that all reviewers agree that the paper is well-written and that side effect avoidance is an important AI safety
4 problem. We are excited that (R1) thinks this work is ground-breaking, novel, and will spur further research, that (R3)
5 is excited by our scaling of AUP, and that our results are considered strong (R1, R5) and significant (R1, R3).

6 **Experiments.** (R1): using multiple auxiliary reward functions performed the same or worse than $|\mathcal{R}| = 1$. (R1, R5)
7 ask for more comparisons. We will include results for DQN and for additional auxiliary reward functions. At (R1)’s
8 suggestion, we tried using the primary reward function as the auxiliary reward function. This condition achieved return
9 and side effects comparable to PPO’s, as its attainable utility shifts did not correlate with side effects.

10 (R5) asks why we didn’t compare to state reachability preservation [14] or to AUP with uniformly randomly drawn
11 auxiliary reward functions over observations (like [25]). Unfortunately, neither approach is remotely viable in SafeLife.
12 We estimate that there are billions of reachable states in any given SafeLife level. Accordingly, we’re aware of a team
13 trying to train reachability-preserving agents, but even the `append-still-easy` task was far too hard. If we generated
14 reward functions by uniformly randomly drawing a reward for every state, the corresponding Q-functions would have
15 extremely high sample complexity in an environment like SafeLife. We used a VAE because the encoder provided
16 sufficient structure for quickly learning the auxiliary Q-function.

17 (R2, R4, R5) ask what distinguishes our work from [25]. The original AUP paper suggested that state reachability
18 preservation avoids the same breadth of side effects as AUP in [25]’s toy environments. We show that unlike reachability
19 preservation, AUP scales to SafeLife; as acknowledged by (R1, R3), we are the first to demonstrate compelling results
20 on any complex side-effect avoidance environment. Furthermore, we showed that AUP competitively accrues reward
21 while avoiding side effects, while equipped with a single auxiliary reward function which was learned *unsupervised*,
22 whereas [25] drew several dozen auxiliary reward functions from the uniform distribution.

23 (R3, R4) wonder about scaling AUP to even more complex environments. We share their interest in this prospect.
24 Realistic settings might have too many side effect opportunities for a supervised penalty to work well. We believe that
25 the efficacy of AUP’s unsupervised penalty term bodes well for even more challenging domains. (R3): we will add
26 more discussion of the challenges AUP may face when scaling further, such as our assumption of a no-op action.

27 (R4) asks what the side effect score (defined lines 115-117) means. Roughly, if AUP halves PPO’s side effect score, then
28 AUP disturbed half as many green cells. Disturbing a patch of green cells corresponds to about 4 additional side effect
29 score. For a qualitative demonstration of the significance of a 46% side effect score reduction, we refer (R4) to the
30 attached GIF files for `append-spawn`.

31 **Training details.** (R1): agents were evaluated on their N_{env} training environments. AUP had no data advantage – it
32 trained for the standard 5M total steps. While AUP_{proj} skips the 100K VAE exploration steps, it still learns its auxiliary
33 Q-function for the first 1M steps. Given more training time, AUP does not overfit and increase side effects. To the
34 contrary, when running three seeds from 5M steps (the paper’s time limit) to 15M steps, AUP’s side effect score changed
35 as follows: `append-still-easy`: -24% , `append-still`: -53% , `append-spawn`: $+17\%$, `prune-still-easy`:
36 -12% . We performed hyperparameter search for `append-still-easy`, and then applied the method successfully to
37 the other three tasks without further tuning. “Random exploration” refers to a uniformly random exploration policy.

38 (R2): the hyperparameter sweep shows the average side effect score and episodic return for the *last* episode, while
39 Figure 3’s charts show a rolling average. The discrepancy was unintentional; we agree it is confusing and will fix it.

40 (R5): we used curriculum training (following the original SafeLife paper, [27]) because PPO fails to learn with just
41 one environment. “Learning R_{AUP} ” refers to the process of training the agent with respect to the AUP reward function
42 defined in equation 1. Note that when training with respect to R_{AUP} , the auxiliary Q-functions are fixed, as they have
43 already been trained. We will clarify the difference between auxiliary and AUP training in camera-ready.

44 **Novelty.** (R2) finds our work lacks novelty because we “take an existing technique and apply it to an existing
45 environment”. This misses our main contribution: demonstrating AUP’s scalability, which is a crucial consideration in
46 AI safety. As acknowledged by (R1, R3, R5), our work provides the *first* strong empirical evidence that AUP (and side
47 effect measures more generally) can scale to any kind of complex environment.

48 **Clarity.** The prior work section and the SafeLife task explanations will be improved until they are as clear as the rest
49 of the paper, and we will incorporate (R1, R5)’s suggested citations. (R2): we will make all code available. We are
50 committed to providing the clearest possible camera-ready paper and look forward to refining the paper to (R1, R2, R3,
51 R4, R5)’s satisfaction.