

Mitigating Goal Misgeneralization via Minimax Regret

Anonymous authors

Paper under double-blind review

Keywords: Goal Misgeneralization, Unsupervised Environment Design, AI Safety, AI Alignment.

Summary

Safe generalization in reinforcement learning requires not only that a learned policy *acts capably* in new situations, but also that it uses its capabilities *towards the pursuit of the designer’s intended goal*. The latter requirement may fail when a *proxy goal* incentivizes similar behavior to the intended goal within the training environment, but not in novel deployment environments. In this setting, policies may behave as if in pursuit of the proxy goal in deployment—a phenomenon known as *goal misgeneralization*. In this paper, we theoretically investigate the possibility of goal misgeneralization under *maximum expected value (MEV)* and *minimax expected regret (MMER)* objectives, and empirically validate our results. Our findings underscore minimax expected regret as a promising principle for mitigating goal misgeneralization.

Contribution(s)

1. We introduce a problem setting called a *proxy-distinguishing distribution shift*, capturing a class of situations in which goal misgeneralization can be elicited and studied.

Context: Given two reward functions (a true goal and a proxy goal), a proxy-distinguishing distribution shift involves (1) training primarily in contexts for which optimizing the proxy goal also optimizes the true goal, then (2) testing primarily in situations in which optimizing the proxy goal is suboptimal under the true goal (so-called *distinguishing levels*). We do not assume training methods have knowledge of the proxy goal.

2. We prove that, under a proxy-distinguishing distribution shift, approximately maximizing expected value on the training distribution permits a misgeneralizing solution if the proportion of distinguishing levels in the training distribution is low enough (Theorem 1).

Context: As a special case, *exactly* maximizing expected value on the training distribution permits a misgeneralizing solution if there are no distinguishing levels in the training distribution. Note: We model *possible* goal misgeneralization—*actual* goal misgeneralization also depends on the agent’s inductive biases for different approximate solutions.

3. We prove that, under a proxy-distinguishing distribution shift, no approximate solution of the minimax expected return objective exhibits goal misgeneralization (Theorem 2).

Context: Theorem 2 holds for fully observable environments; we include a generalization to partially observable environments in the supplementary materials (Theorem 3).

4. We show empirically that, under conditions approximating a proxy-distinguishing distribution shift in procedurally generated grid-world environments, policies learned using domain randomization (DR; an MEV-based training method) exhibit goal misgeneralization when the proportion of distinguishing levels in the training distribution is low enough (§7.1).

Context: Langosco et al. (2022) demonstrated goal misgeneralization with zero distinguishing levels, we extend this finding to the case with a small positive proportion.

5. We show empirically that, under the same conditions, existing regret-based unsupervised environment design (UED) methods, PLR^\perp (Jiang et al., 2022) and ACCEL (Parker-Holder et al., 2023), (1) can detect rare distinguishing levels and increase their proportion in the training distribution, and (2) are more robust to goal misgeneralization than DR is (§7.2).

Context: In some cases, less advanced UED methods fail to find MMER policies, and still exhibit goal misgeneralization (§7.3, §7.4), indicating that more mature UED methods are needed to achieve the potential of MMER for preventing goal misgeneralization in practice.

Mitigating Goal Misgeneralization via Minimax Regret

Anonymous authors

Paper under double-blind review

Abstract

Safe generalization in reinforcement learning requires not only that a learned policy acts capably in new situations, but also that it uses its capabilities towards the pursuit of the designer’s *intended goal*. The latter requirement may fail when a *proxy goal* incentivizes similar behavior to the intended goal within the training environment, but not in novel deployment environments. This creates the risk that policies will behave as if in pursuit of the proxy goal, rather than the intended goal, in deployment—a phenomenon known as *goal misgeneralization*. In this paper, we formalize this problem setting in order to theoretically study the possibility of goal misgeneralization under different training objectives. We show that goal misgeneralization is possible under approximate optimization of the maximum expected value (MEV) objective, but not the minimax expected regret (MMER) objective. We then empirically show that the standard MEV-based training method of domain randomization exhibits goal misgeneralization in procedurally-generated grid-world environments, whereas current regret-based unsupervised environment design (UED) methods are more robust to goal misgeneralization (though they don’t find MMER policies in all cases). Our findings suggest that minimax expected regret is a promising approach to mitigating goal misgeneralization.

1 Introduction

As reinforcement learning (RL) is increasingly applied in complex, open-ended, real-world environments, it is becoming infeasible for training to comprehensively cover all situations an agent will face in deployment. We therefore need training methods to produce policies that *generalize*, behaving as intended when faced with a novel scenario (Kirk et al., 2023).

A particular challenge arises when incomplete coverage of the environment space during training creates a *proxy goal*. A proxy goal is a reward function that, compared to the true goal, induces similar optimal behavior in most situations encountered during training, but induces radically different behavior in some novel situations. Proxy goals create the risk of *goal misgeneralization*—learning a policy that retains its capabilities in novel situations, but behaves as if to pursue the proxy goal instead of the true goal when the two diverge (Langosco et al., 2022; Shah et al., 2022).

Goal misgeneralization can arise when such “goal-distinguishing” situations—where the proxy goal and the true goal diverge—are rare within training, making policies that pursue the wrong goal approximately optimal in terms of the standard RL objective of maximum expected value (MEV). This motivates the need for training methods that can somehow identify goal-distinguishing situations within a complex environment, and ensure they are adequately covered in the training distribution.

We observe that favoring the proxy goal in goal-distinguishing situations leads to high *expected regret*, defined as the shortfall of expected return compared to that obtained by an optimal policy. An environment selected to maximize a policy’s expected regret will naturally include goal-distinguishing situations as long as the policy ignores the true goal. Therefore, we propose mitigating goal misgeneralization via the *minimax expected regret* (MMER; Savage, 1951) objective.

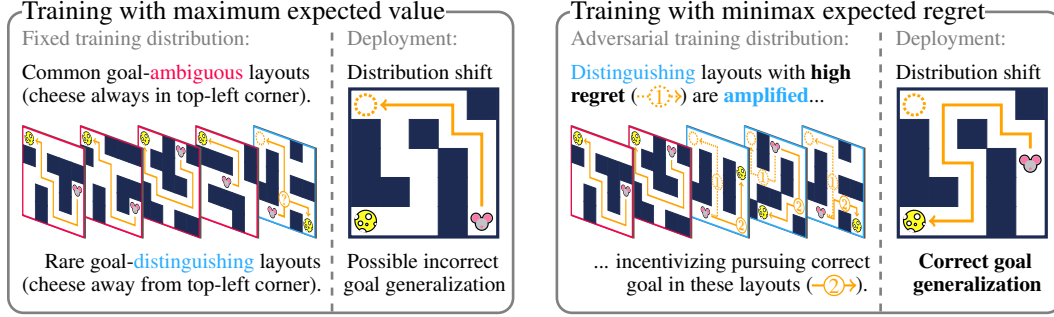


Figure 1: **Maximum expected value and minimax expected regret vs. goal misgeneralization.** A mouse searches a maze for cheese that is usually located in the top-left corner. There is a proxy goal (“go to the corner”) that mostly incentivizes the same optimal behavior as the true goal (“go to the cheese”). (Left): Standard RL methods that *approximately* maximize expected value/return could find a policy that behaves as if pursuing the proxy goal rather than the true goal, since layouts where this policy fails are rare in training. This would lead to incorrect generalization. (Right): If a policy ignores the cheese, a regret-maximizing adversary can move the cheese away from the corner until the agent internalizes the correct goal, leading to correct generalization.

- 38 In this paper, we conduct a theoretical and empirical investigation of the possibility of goal misgen-
 39 eralization under the MEV and MMER objectives. An outline of our contributions is as follows.
- 40 1. In Section 4, we introduce a problem setting called a *proxy-distinguishing distribution shift*,
 41 formalizing a class of situations in which goal misgeneralization can arise.
 - 42 2. In Section 5, we show formally that (1) approximately optimizing MEV is susceptible to goal
 43 misgeneralization if goal-distinguishing situations are sufficiently rare (Theorem 1), and (2) ap-
 44 proximately optimizing MMER is provably robust to goal misgeneralization (Theorem 2).
 - 45 3. In Sections 6 and 7, we empirically study the robustness to goal misgeneralization of a standard
 46 MEV-based training method (domain randomization; Tobin et al., 2017), and recent MMER-
 47 based training methods (regret-based unsupervised environment design; Dennis et al., 2020; Jiang
 48 et al., 2022; Parker-Holder et al., 2023) under a proxy-distinguishing distribution shift.
- 49 Our theoretical results show that, in the limit of idealized training methods, MMER-based train-
 50 ing is guaranteed to be robust against goal misgeneralization, whereas MEV-based training is not.
 51 Our empirical results show that current MMER-based training methods are indeed more robust to
 52 goal misgeneralization than MEV-based training is, and, while they sometimes still exhibit goal
 53 misgeneralization, this happens less for more advanced methods. Together, these results establish
 54 MMER-based training as a promising approach to preventing goal misgeneralization.

55 2 Related work

56 **Goal misgeneralization.** Ensuring learned systems generalize as intended in novel situations is a
 57 perennial challenge for deep learning and deep RL (Kirk et al., 2023). Christiano (2018) distin-
 58 guishes between *benign* generalization failures, where an agent fails to behave capably in a novel
 59 situation, and *malign* generalization failures, where the agent demonstrates capable behavior towards
 60 the pursuit of an unintended objective. Langosco et al. (2022) and Shah et al. (2022) demonstrate
 61 behavioral examples of malign generalization failures in deep RL, introducing the term *goal misgen-*
 62 *eralization*. Goal misgeneralization is similar to *shortcut learning* in supervised learning (Geirhos
 63 et al., 2020), but emphasizes shortcut *reward functions*, rather than shortcut policies.

64 Recent work proposes complementary approaches to mitigating the risk of goal misgeneralization.
 65 Starace (2023) investigates influencing the agent’s inductive bias in favor of correct goal generaliza-
 66 tion using goal-conditioned RL with natural language task descriptions. Trinh et al. (2024) studies

67 methods for detecting when the agent is in an unfamiliar situation and choosing to ask an expert (at
68 a cost) to clarify the optimal action.

69 **Training in complex environments.** The standard technique for RL complex environments is to
70 train on situations sampled from a fixed distribution, a technique known as domain randomization
71 (e.g., [Tobin et al., 2017](#); [Peng et al., 2018](#)). Maximizing expected return over such situations corre-
72 sponds to pursuing the MEV objective with respect to the fixed training distribution.

73 [Dennis et al. \(2020\)](#) proposed *regret-based unsupervised environment design (UED)*, an RL training
74 technique featuring an adversarial environment designer that continually adapts the training distribu-
75 tion aiming to maximize the agent’s expected regret. Maximizing expected return on this adversarial
76 distribution corresponds to the MMER objective ([Dennis et al., 2020](#)). UED has been promoted as
77 a technique for (1) improving sample efficiency by creating an emergent curriculum; and (2) im-
78 proving capability generalization via adversarial robustness ([Dennis et al., 2020](#); [Jiang et al., 2022](#);
79 [Parker-Holder et al., 2023](#)). We show that UED also helps to mitigate goal misgeneralization.

80 Alternative adversarial approaches, such as maximin expected value ([Dennis et al., 2020](#); [Wang](#)
81 [et al., 2023](#)), maximizing diversity ([OpenAI et al., 2019](#)), or maximizing learnability ([Rutherford](#)
82 [et al., 2024](#)), have not been studied in the context of goal misgeneralization. These approaches may
83 also mitigate goal misgeneralization to the extent that they promote training in goal-distinguishing
84 situations, incentivizing the agent to internalize the true goal. We show that directly optimizing the
85 training environment for regret is sufficient. Appendix I shows that minimax expected value can
86 exhibit goal misgeneralization when some situations have low maximum expected return.

87 3 Preliminaries

88 A (reward-free) underspecified Markov decision processes (UMDP) is a tuple $M = \langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$
89 where Θ is a space of free parameters (also called **levels**), \mathcal{A} is the agent’s action space, \mathcal{S} is a
90 state space, $\mathcal{I} : \Theta \rightarrow \Delta(\mathcal{S})$ is an initial state distribution, and $\mathcal{T} : \Theta \times \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ is a
91 conditional transition distribution. Given a level $\theta \in \Theta$ we have a fully-specified (reward-free)
92 MDP $\langle \mathcal{A}, \mathcal{S}, \mathcal{I}(\theta), \mathcal{T}(\theta, -, -) \rangle$. We aggregate these MDPs into a single complex environment using
93 a **level distribution** $\Lambda \in \Delta(\Theta)$. A **reward function** (or **goal**) is a function $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$.
94 Taken together, M and R define a proper (non-reward-free) UMDP. We define reward functions
95 and reward-free UMDPs separately to facilitate considering multiple goals for an otherwise fixed
96 environment. We usually denote by R and \tilde{R} the **true goal** and the **proxy goal**, respectively.

97 An agent’s **policy** is a conditional action distribution $\pi : \Theta \times \mathcal{S} \rightarrow \Delta(\mathcal{A})$. Note that we assume
98 the policy observes the level (we consider the partially observable case in Appendix C). The set of
99 all policies is denoted by Π . We define the **expected return** (or **expected value**) of policy π in the
100 level θ under goal R as the discounted cumulative reward $V^R(\pi; \theta) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1})]$
101 where $\gamma \in (0, 1)$ is a discount factor and the expectation is over $s_0 \sim \mathcal{I}(\theta)$, $a_t \sim \pi(\theta, s_t)$, and
102 $s_{t+1} \sim \mathcal{T}(\theta, s_t, a_t)$. We lift this definition to level distributions as $V^R(\pi; \Lambda) = \mathbb{E}_{\theta \sim \Lambda}[V^R(\pi; \theta)]$.
103 A **normalized goal** is one such that the return has support in $[0, 1]$.

104 We define the **expected regret** of a policy π in the level θ under a goal R as the shortfall of expected
105 value achieved by the policy compared to an optimal policy for that level,

$$G^R(\pi; \theta) = \max_{\pi' \in \Pi} V^R(\pi'; \theta) - V^R(\pi; \theta). \quad (1)$$

106 Once again, we lift this definition to level distributions as $G^R(\pi; \Lambda) = \mathbb{E}_{\theta \sim \Lambda}[G^R(\pi; \theta)]$. Since the
107 policy is conditioned on θ , we also have the following identity (see Appendix B):

$$G^R(\pi; \Lambda) = \max_{\pi' \in \Pi} V^R(\pi'; \Lambda) - V^R(\pi; \Lambda). \quad (2)$$

4 Problem setting

Langosco et al. (2022) and Shah et al. (2022) formalize goal misgeneralization and provide case studies in which it arises. In order to theoretically study goal misgeneralization, we formalize an abstract class of situations in which goal misgeneralization can arise as a problem setting called a *proxy-distinguishing distribution shift*.

Given an UMDP, consider a pair of reward functions: one true goal and one proxy goal. We classify levels by whether the proxy goal incentivizes optimal or suboptimal behavior under the true goal.

Definition 1 (Ambiguous level). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a true goal R , and a proxy goal \tilde{R} . A level $\theta \in \Theta$ is (universally, perfectly, goal-) ambiguous with respect to R and \tilde{R} if all optimal policies with respect to the proxy goal are also optimal with respect to the true goal, that is*

$$\arg \max_{\pi \in \Pi} V^{\tilde{R}}(\pi; \theta) \subseteq \arg \max_{\pi \in \Pi} V^R(\pi; \theta).$$

Definition 2 (C -distinguishing level). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a true goal R , a proxy goal \tilde{R} , and a constant $C \geq 0$. A level $\theta \in \Theta$ is (universally goal-) C -distinguishing with respect to R and \tilde{R} if all policies optimal with respect to \tilde{R} achieve C -suboptimal expected return with respect to R , that is,*

$$\forall \pi \in \arg \max_{\pi \in \Pi} V^{\tilde{R}}(\pi; \theta), \quad V^R(\pi; \theta) < \max_{\pi' \in \Pi} V^R(\pi'; \theta) - C.$$

We note that this classification is not exhaustive: for a given UMDP and pair of reward functions, some levels may be neither perfectly ambiguous nor C -distinguishing for any $C \geq 0$; the conditions must hold for all optimal policies (hence “universally”). However, it is often true that optimizing a misspecified goal leads to arbitrarily low return (cf. Zhuang & Hadfield-Menell, 2020).

We model the shift from training to deployment as a change in level distribution. The training distribution represents all levels an RL method can access for training a policy prior to deployment. The distribution shift is *proxy-distinguishing* when the training distribution concentrates mostly on ambiguous levels but the deployment distribution concentrates mostly on distinguishing levels.

Definition 3 (Proxy-distinguishing distribution shift). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a true goal R , and a proxy goal \tilde{R} . A proxy-distinguishing distribution shift is a tuple $\langle \alpha, \beta, C, \Lambda^{\text{Train}}, \Lambda^{\text{Deploy}} \rangle$ where α, β are ratios such that $0 \leq \alpha < \beta \leq 1$, $C \geq 0$ is a constant, and $\Lambda^{\text{Train}}, \Lambda^{\text{Deploy}} \in \Delta(\Theta)$ are level distributions over levels with the following classifications (with respect to R and \tilde{R}):*

1. Λ^{Train} has probability α on C -distinguishing levels and probability $1 - \alpha$ on ambiguous levels.
2. Λ^{Deploy} has probability β on C -distinguishing levels and probability $1 - \beta$ on ambiguous levels.

We are mainly interested in the case where α is very close to zero (where goal misgeneralization is a particular risk) and β is very close to one (where goal misgeneralization is a particular concern).

We don’t assume prior knowledge of the proxy goal or distinguishing levels. However, we *do* assume the ability to train in distinguishing levels, once identified. In practice, one can train in a very wide space of situations, whether via a simulator (Tobin et al., 2017; Peng et al., 2018; Kumar et al., 2021; Makoviychuk et al., 2021; Muratore et al., 2022; Ma et al., 2024), a generative environment model (Bruce et al., 2024), or a world model (Ha & Schmidhuber, 2018; Hafner et al., 2019; Schrittwieser et al., 2020; Hafner et al., 2023; Valevski et al., 2024). If *all* levels accessible before deployment are ambiguous, we may require alternative assumptions (see, e.g., Trinh et al., 2024).

Moreover, we assume access to a reliable reward signal in favor of the true goal in distinguishing levels. This mirrors assumptions made in work on spurious correlations in supervised learning (e.g., Liu et al., 2021; Zhang et al., 2022). However, in practice, reward functions may be subject to misspecification in such corner cases (cf. Hadfield-Menell et al., 2017). Future work could develop methods that treat the true goal as underspecified in rare, distinguishing levels and find ways to incentivize safe generalization behavior despite this uncertainty.

5 Theoretical results

In this section, we prove that under a proxy-distinguishing distribution shift, the maximum expected value (MEV) objective permits an approximately optimal policy that exhibits goal misgeneralization. On the other hand, we show that any policy that is approximately optimal with respect to minimax expected regret (MMER) must avoid goal misgeneralization. All proofs are in Appendix A.

We consider *approximately* optimal policies because, in practice, training uses finite optimization power and will not always find policies that are exactly optimal according to the given objective. We model approximate optimization by supposing that our optimization methods will instead find an arbitrary policy within a small threshold of optimal for the given objective. We use the notation $\arg\text{-}\varepsilon\text{-}\max_{x \in \mathcal{X}} f(x) = \{x \in \mathcal{X} \mid f(x) \geq \max_{\xi \in \mathcal{X}} f(\xi) - \varepsilon\}$ (likewise $\arg\text{-}\varepsilon\text{-}\min$) for approximate optimization of a function $f : \mathcal{X} \rightarrow \mathbb{R}$ with approximation threshold $\varepsilon \geq 0$.

5.1 Approximate maximum expected value is susceptible to goal misgeneralization

The standard objective used in RL is maximum expected value (MEV) with respect to the fixed training distribution. We formalize this objective as follows.

Definition 4 (Approximate MEV). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a goal R , an approximation threshold $\varepsilon \geq 0$, and a fixed level distribution $\Lambda \in \Delta(\Theta)$. The MEV objective with respect to Λ is to maximize $V^R(\pi; \Lambda)$. The approximate MEV policy set with respect to Λ is then*

$$\Pi_{\varepsilon}^{\text{MEV}}(R, \Lambda) = \arg\text{-}\varepsilon\text{-}\max_{\pi \in \Pi} V^R(\pi; \Lambda).$$

The MEV objective permits goal misgeneralization under a proxy-distinguishing distribution shift if the proportion of distinguishing levels in training is too small. Intuitively, a policy that pursues the proxy goal in all levels achieves enough return on ambiguous levels to be approximately optimal. Note: rather than modeling inductive bias, we characterize the *possibility* of goal misgeneralization.

Theorem 1 (MEV is susceptible to goal misgeneralization). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a pair of normalized goals R, \tilde{R} , a proxy-distinguishing distribution shift $\langle \alpha, \beta, C, \Lambda^{\text{Train}}, \Lambda^{\text{Deploy}} \rangle$, and an approximation threshold $\varepsilon \geq 0$. If $\varepsilon \geq \alpha$, then*

$$\exists \pi^{\text{MEV}} \in \Pi_{\varepsilon}^{\text{MEV}}(R, \Lambda^{\text{Train}}) \text{ such that } V^R(\pi^{\text{MEV}}; \Lambda^{\text{Deploy}}) < \max_{\pi \in \Pi} V^R(\pi; \Lambda^{\text{Deploy}}) - \beta C.$$

5.2 Approximate minimax expected regret is robust to goal misgeneralization

The minimax expected regret (MMER) objective can be defined via an equilibrium of a two-player zero-sum game, in which an agent selects a policy and an adversary selects a level distribution. We relax this definition by allowing both players to play an *approximate* best response.

Definition 5 (Approximate MMER). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a goal R , and approximation thresholds $\varepsilon, \delta \geq 0$. Consider the two-player zero-sum game $\langle \langle \Pi, \Delta(\Theta) \rangle, \langle -G^R, G^R \rangle \rangle$, where an agent plays a policy $\pi \in \Pi$ and an adversary plays a level distribution $\Lambda \in \Delta(\Theta)$, aiming to minimize or maximize $G^R(\pi; \Lambda)$ respectively. A pair (π, Λ) is an (ε, δ) -equilibrium if $\pi \in \arg\text{-}\varepsilon\text{-}\min_{\pi' \in \Pi} G^R(\pi'; \Lambda)$ and $\Lambda \in \arg\text{-}\delta\text{-}\max_{\Lambda' \in \Delta(\Theta)} G^R(\pi; \Lambda')$. The approximate MMER policy set is then*

$$\Pi_{\varepsilon, \delta}^{\text{MMER}}(R) = \{\pi \in \Pi \mid \exists \Lambda \in \Delta(\Theta) \text{ such that } (\pi, \Lambda) \text{ is an } (\varepsilon, \delta)\text{-equilibrium}\}.$$

The MMER objective does not permit goal misgeneralization under any distribution shift within the adversary's strategy space. Intuitively, if the agent's policy pursues the proxy goal on distinguishing levels, the adversary can exploit the high regret until the agent's policy pursues the true goal.

Theorem 2 (MMER is robust to goal misgeneralization). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a pair of goals R, \tilde{R} , a proxy-distinguishing distribution shift $\langle \alpha, \beta, C, \Lambda^{\text{Train}}, \Lambda^{\text{Deploy}} \rangle$, and approximation thresholds $\varepsilon, \delta \geq 0$. Then*

$$\forall \pi^{\text{MMER}} \in \Pi_{\varepsilon, \delta}^{\text{MMER}}(R), \text{ we have } V^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) \geq \max_{\pi \in \Pi} V^R(\pi; \Lambda^{\text{Deploy}}) - \varepsilon - \delta.$$

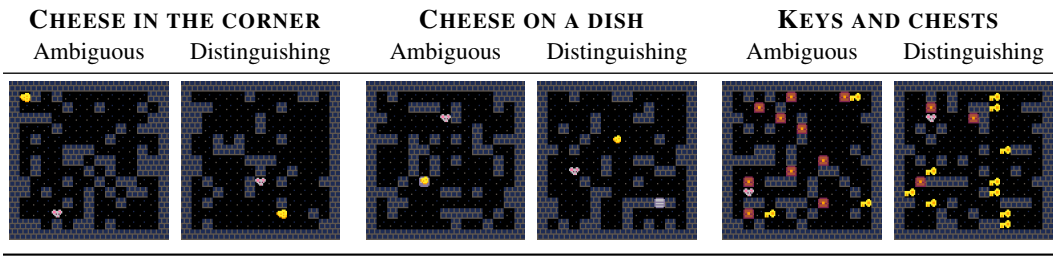


Figure 2: Example procedurally-generated ambiguous/distinguishing levels. The agent’s observation is a $15 \times 15 \times c$ Boolean grid (where c is an environment-dependent number of channels).

192 A slightly modified bound holds for partially observable environments after accounting for the min-
 193 imum expected regret *realizable by a fixed policy* (see Appendix C).

194 6 Experimental methods

195 In this section, we outline our methods for investigating the robustness to goal misgeneralization of
 196 MEV-based and MMER-based training methods. We construct three custom procedurally-generated
 197 grid-world environments approximating proxy-distinguishing distribution shifts (Section 6.1). We
 198 compare a standard MEV-based training method and two recently proposed MMER-based methods
 199 with adversaries of varying flexibility (Section 6.2), paired with regret estimators that leverage vary-
 200 ing amounts of domain knowledge (either using the ground truth maximum return, or estimating it
 201 from samples; Section 6.3). Section 7 presents the results of our experiments.

202 6.1 Procedurally-generated grid-world environments

203 Langosco et al. (2022) exhibited goal misgeneralization in several environments from OpenAI Proc-
 204 gen (Cobbe et al., 2020), suitably modified to implement a proxy-distinguishing distribution shift
 205 with $\alpha = 0$. We implement three similar procedurally-generated grid-world environments in JAX
 206 (Bradbury et al., 2018), allowing us to more easily implement custom level generation and analysis.

207 For each environment, we construct two procedural level generators $\Lambda_{\text{Ambig.}}, \Lambda_{\text{Distg.}} \in \Delta(\Theta)$, ap-
 208 proximately concentrated on ambiguous and distinguishing levels, respectively. From these, we
 209 define training distributions $\Lambda_{\alpha}^{\text{Train}} = (1 - \alpha)\Lambda_{\text{Ambig.}} + \alpha\Lambda_{\text{Distg.}}$, where α is the proportion of distin-
 210 guishing levels. In our experiments, we vary α between 10^{-5} to 10^{-1} , with $\alpha \in \{0, 1\}$ as baselines.
 211 We evaluate on $\Lambda^{\text{Deploy}} = \Lambda_{\text{Distg.}}$, approximating a proxy-distinguishing distribution shift.

212 The three environments are as follows. Figure 2 illustrates example levels (note we use Boolean
 213 observations). Appendix D comprehensively documents each environment, including the details of
 214 classifying levels as ambiguous or distinguishing and procedural level generation.

- 215 1. **CHEESE IN THE CORNER.** A mouse navigates a maze. The true goal assigns +1 reward for
 216 reaching a piece of cheese, while a proxy goal assigns +1 reward for reaching the top left corner
 217 for the first time. Levels with the cheese in the top left corner are ambiguous and levels with the
 218 cheese away from the corner are mostly distinguishing.
- 219 2. **CHEESE ON A DISH.** This time the mouse navigates a maze containing cheese and also a dish.
 220 The true goal assigns +1 reward for reaching the cheese, while a proxy goal assigns +1 reward
 221 for reaching the dish. Levels with the cheese and dish co-located are ambiguous, and levels with
 222 the cheese and dish separated are mostly distinguishing.
- 223 3. **KEYS AND CHESTS.** A more complex, multi-stage task, in which the mouse navigates a maze
 224 collecting keys and spending them to open chests. Levels with 3 keys and 10 chests are approxi-
 225 mately ambiguous. Levels with 10 keys and 3 chests are mostly distinguishing—a misgeneraliz-
 226 ing policy would overprioritize key collection beyond what is necessary for opening chests.

6.2 Training methods

For both MEV-based and MMER-based training, we follow [Langosco et al. \(2022\)](#) and use an agent network architecture based on that of IMPALA ([Espeholt et al., 2018](#)) with a dense feed-forward layer replacing the LSTM block. We perform policy updates with PPO ([Schulman et al., 2017](#)) and GAE ([Schulman et al., 2015](#)). We document hyperparameters and compute usage in Appendix E.

For MEV, we use a standard method for training in UMDPs given a fixed level distribution.

1. **Domain randomization** (DR; [Tobin et al., 2017](#)). For each iteration of PPO, we sample (procedurally generate) a new batch of levels from the fixed training level distribution $\Lambda_{\alpha}^{\text{Train}}$, collect experience in this batch of levels, and then train on the collected experience.

For MMER, we use two methods of regret-based unsupervised environment design (UED; [Dennis et al., 2020](#)). UED methods implement the two-player zero-sum game from Definition 5 by training the policy on levels selected from a distribution chosen by a regret-maximizing **adversary**. The first UED method is a regret-based form of prioritized level replay ([Jiang et al., 2021](#)).

2. **Robust prioritized level replay** (PLR⁺; [Jiang et al., 2022](#)). The adversary parametrizes its level distribution using a fixed-size **level buffer**. Throughout training, the adversary refines the buffer by either (1) sampling a new batch of levels from the underlying training distribution $\Lambda_{\alpha}^{\text{Train}}$ and estimating the expected regret of the current policy on these levels; or (2) sampling from the current buffer, conducting a PPO training step with the chosen levels, and updating their expected regret estimates; keeping the highest-regret levels in the buffer.

PLR⁺ has the advantage of being domain-agnostic, but has the disadvantage of only being able to *replay* levels once they have been sampled from the underlying distribution. We also consider a more advanced adversary with an independent means of exploring the space of level distributions.

3. **Adversarially compounding complexity by editing levels** (ACCEL; [Parker-Holder et al., 2023](#)). The adversary continually refines a level buffer with steps (1) and (2) from PLR⁺, and additionally by (3) applying stochastic **edits** to the levels used for PPO training to generate similar levels, and estimating the expected regret of the current policy on these new levels.

ACCEL additionally requires an **edit distribution**. We edit levels by sampling a sequence of random elementary level modifications, none of which change whether the level is ambiguous or distinguishing. Appendix H details this edit distribution and compares it to edit distributions with more or less ability to introduce distinguishing levels.

6.3 Expected regret estimation methods

Both UED methods require an (**expected**) **regret estimator** for deciding which levels to keep in the buffer. To represent the current capabilities of UED methods, we use the following domain-agnostic estimator, similar to the MaxMC estimator proposed by [Jiang et al. \(2022\)](#).

1. **Max-latest estimator**. We estimate the expected regret of policy π in level θ under goal R as

$$\hat{G}_{\text{max-latest}}^R(\pi; \theta) = \hat{V}_{\text{max}}^R(\theta) - \hat{V}_{\text{latest}}^R(\pi; \theta) \quad (3)$$

where $\hat{V}_{\text{max}}^R(\theta)$ is the highest empirical return ever achieved for this level throughout training; and $\hat{V}_{\text{latest}}^R(\pi; \theta)$ is the empirical average return achieved by the current policy.

To simulate a more advanced regret estimator than is currently available in practice, we also consider a domain-specific estimator that solves each procedurally-generated level using a graph algorithm to compute the exact maximum expected return (details in Appendix D).

2. **Oracle-latest estimator**. We estimate the expected regret of policy π in level θ under goal R as

$$\hat{G}_{\text{oracle-latest}}^R(\pi; \theta) = \max_{\pi'} V^R(\pi'; \theta) - \hat{V}_{\text{latest}}^R(\pi; \theta) \quad (4)$$

where $\max_{\pi'} V^R(\pi'; \theta)$ is the maximum expected return for the level; and $\hat{V}_{\text{latest}}^R(\pi; \theta)$ is as above.

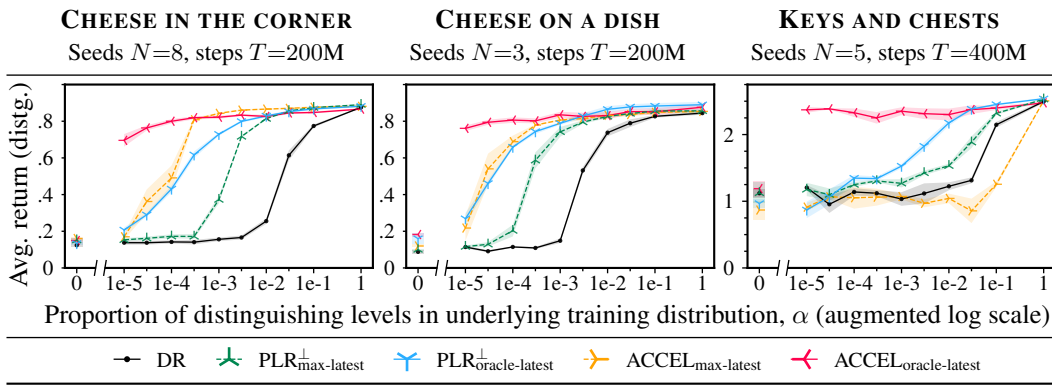


Figure 3: **Distribution shift performance for various training distributions.** Average return over 512 steps for an evaluation batch of 256 distinguishing levels sampled from $\Lambda^{\text{Deploy}} = \Lambda^{\text{Distg.}}$. High performance indicates policies generalizing as intended; low performance indicates goal misgeneralization. Each policy is trained on T environment steps using the indicated training method with underlying training distribution $\Lambda_{\alpha}^{\text{Train}} = (1 - \alpha)\Lambda^{\text{Ambig.}} + \alpha\Lambda^{\text{Distg.}}$. Mean over N seeds, shaded to one standard error. Note the split in the horizontal axis used to show zero on the log scale.

269 7 Experimental results

270 In this section, we report the results of our main experiments. Consistent with Theorem 1, MEV-
 271 based training is susceptible to goal misgeneralization unless the proportion of distinguishing levels
 272 in the training distribution is sufficiently high (Section 7.1). Consistent with Theorem 2, MMER-
 273 based training methods are typically capable of identifying and increasing the proportion of rare,
 274 high-regret, distinguishing levels, thereby preventing goal misgeneralization in many situations
 275 where MEV-based training misgeneralizes (Section 7.2). In some cases, UED methods fail to find
 276 MMER policies, and exhibit goal misgeneralization. We see generally that the more advanced UED
 277 methods are more robust to goal misgeneralization (Section 7.3). In KEYS AND CHESTS, MEV
 278 outperforms ACCEL with max-latest regret estimation, underscoring reliable regret estimation as a
 279 particular challenge for future work on MMER-based training (Section 7.4).

280 7.1 Domain randomization exhibits goal misgeneralization with rare distinguishing levels

281 Theorem 1 says that if the proportion of distinguishing levels in the fixed training distribution is
 282 small enough, then approximately optimizing MEV *possibly* leads to goal misgeneralization. Our
 283 experiments show that DR, an MEV-based training method, indeed exhibits goal misgeneralization
 284 when the proportion of distinguishing levels in the training distribution is small enough. Figure 3
 285 shows end-of-training performance on distinguishing levels. There is a threshold below which DR
 286 performance on distinguishing levels falls. DR achieves high return on ambiguous levels and high
 287 proxy return on distinguishing levels (see Appendix F), indicating a case of goal misgeneralization.

288 In CHEESE IN THE CORNER and KEYS AND CHESTS, DR exhibits goal misgeneralization until there
 289 is around $\alpha = 1e-1$ (10%) mass on distinguishing levels. For CHEESE ON A DISH, DR is robust to
 290 goal misgeneralization from as low as $\alpha = 1e-2$ (1%) (see also Appendix L). Appendix J shows that
 291 training for substantially longer slightly increases DR’s robustness in CHEESE ON A DISH.

292 We note that Langosco et al. (2022) previously demonstrated goal misgeneralization while training
 293 without distinguishing levels in similar environments. Moreover, Langosco et al. (2022) demon-
 294 strated that for a modified version of OpenAI ProcGen’s COINRUN environment (Cobbe et al., 2019;
 295 2020), training with $\alpha = 2e-2$ (2%) prevents goal misgeneralization. They did not experiment with
 296 smaller proportions of distinguishing levels. We show that with small but nonzero proportions of
 297 distinguishing levels, domain randomization can still exhibit goal misgeneralization.

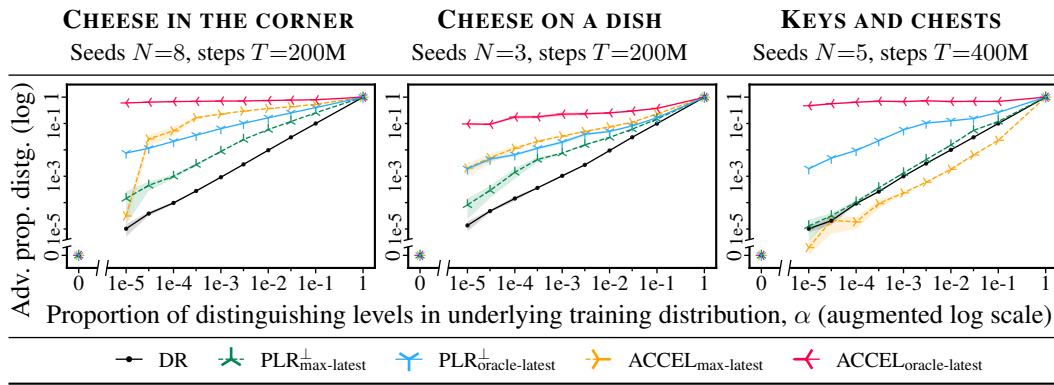


Figure 4: **Rate at which distinguishing levels are played by adversary.** We classify levels played by the adversary as either ambiguous or distinguishing. We plot the proportion of distinguishing levels across training for T environment steps. The diagonal represents the proportion of distinguishing levels sampled from the underlying training distribution $\Lambda_{\alpha}^{\text{Train}} = (1 - \alpha)\Lambda_{\text{Ambig.}} + \alpha\Lambda_{\text{Distg.}}$ (these levels are used for training in DR). Points above the diagonal indicate the adversary increasing the proportion of distinguishing levels relative to the underlying training distribution. Mean over N seeds, shaded to one standard error. Note the splits in *both* axes used to show zero on the log scales.

7.2 Regret-based prioritization amplifies distinguishing levels, mitigating misgeneralization

Theorem 2 says that, at an approximate equilibrium of the MMER game, the agent will play a policy that always pursues the true goal. Otherwise, the adversary could play a distribution of distinguishing levels, leading to high regret. Figure 4 shows the average proportion of distinguishing levels played by the adversary throughout training, showing that, with the exception of max-latest estimation in the KEYS AND CHESTS environment, the adversary plays distinguishing levels disproportionately often compared to sampling from the underlying distribution.

Figure 3 shows that this increase in the proportion of training levels is, in most cases, enough to lead to a policy that pursues the intended goal. In each environment, MMER-based training methods are robust to goal misgeneralization at α values for which DR exhibits goal misgeneralization. For example, in CHEESE IN THE CORNER, all UED methods are robust to goal misgeneralization at $\alpha = 1e-2$ (1%), and some remain robust for even lower α . Note that some evaluation levels are unsolvable—the highest return to be expected is given by the agents trained with $\alpha = 1$.

7.3 Increasingly advanced UED methods are more robust to goal misgeneralization

Theorem 2 says that MMER-based training should be robust to goal misgeneralization regardless of the distribution shift. In contrast, in our experiments, the proportion of distinguishing levels played by the adversary decreases as we decrease α (Figure 4), and each UED method exhibits a threshold below which it fails to converge to an MMER policy, and exhibits goal misgeneralization (Figure 3).

This performance trend reflects how the adversaries construct level distributions. When distinguishing levels are very rare (or never arise), the adversary is hindered (prevented) from increasing the number of distinguishing levels in the buffer. Compared to PLR^{\perp} , ACCEL can replicate similar levels throughout its buffer through edits, but we used edits that don't create new distinguishing levels. Appendix H investigates ACCEL variants with different edit distributions, showing that ACCEL can prevent goal misgeneralization even when $\alpha = 0$ if edits can introduce distinguishing levels.

Overall, robustness correlates with how advanced the adversaries are. The most flexible adversary (ACCEL) paired with the most powerful expected regret estimator (oracle-latest) is remarkably robust to goal misgeneralization in all environments for all positive α tested. The less flexible PLR^{\perp} using the less powerful max-latest expected regret estimator is the least robust.

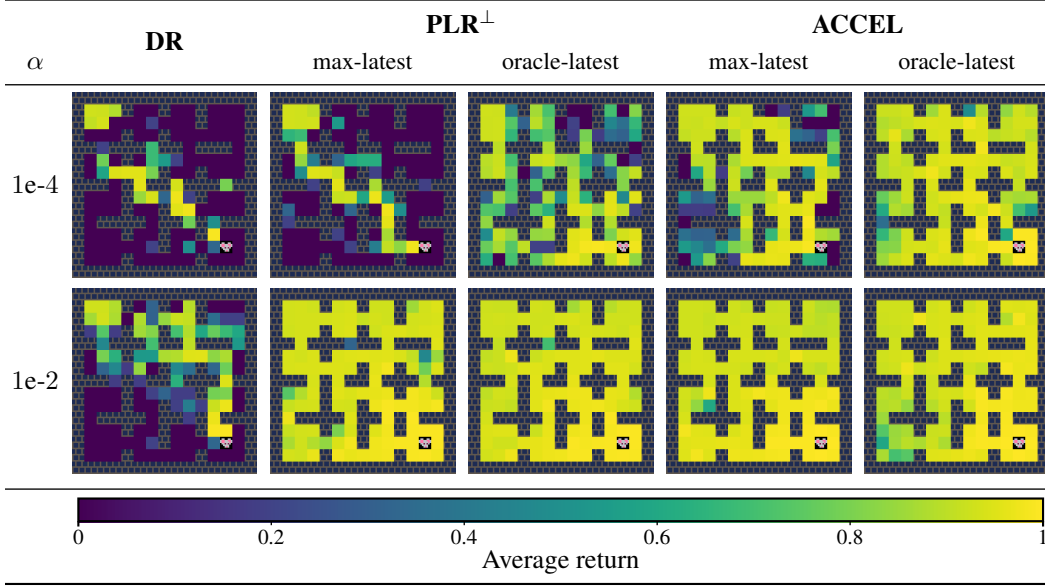


Figure 5: **Performance on CHEESE IN THE CORNER levels with varying cheese position.** For each training configuration, we evaluate the trained policy (first of 8 seeds) on a batch of 122 levels with shared wall layout and mouse spawn position but different cheese positions. We indicate average return on levels with different cheese positions by the color of the corresponding grid square. We see a progression whereby for more advanced algorithms or higher α , the agent is robust to a greater proportion of cheese positions. See Appendix G for more details and full range of α values.

326 7.4 Biased regret estimation can undermine UED in more complex environments

327 The poor performance of ACCEL with max-latest the estimator in the KEYS AND CHESTS environ-
 328 ment underscores the challenge of expected regret estimation. Estimating maximum return from
 329 samples is particularly challenging in this environment, where high return is unlikely to be achieved
 330 in distinguishing levels by chance, since chests are substantially rarer than in ambiguous levels. It
 331 appears that the increased flexibility of ACCEL in this case works as a disadvantage, leading to the
 332 adversary being led astray by biased regret estimates even more so than PLR^\perp .

333 The challenges of regret estimation are known, and are an active area of research (cf. [Rutherford](#)
 334 [et al., 2024](#)). Our results highlight the importance of future work on reliable regret estimation meth-
 335 ods, towards achieving the improved performance shown by our domain-specific oracle-latest esti-
 336 mator. Such work could investigate using a separate policy network to estimate the maximum return
 337 (cf. [Dennis et al., 2020](#)), or incorporating the predictions of a value network (cf. [Jiang et al., 2022](#)).

338 8 Conclusion

339 In this paper, we introduce the setting of a proxy-distinguishing distribution shift, and offer a the-
 340oretical and empirical investigation of the robustness of MEV-based and MMER-based training to
 341 goal misgeneralization. We show theoretically and empirically that MEV-based training on a fixed
 342 training distribution can lead to goal misgeneralization. In contrast, we show that MMER-based
 343 training is provably robust against goal misgeneralization in the limit of idealized training meth-
 344 ods, and regret-based unsupervised environment design methods are empirically more robust than
 345 MEV-based training. Current UED methods do not find MMER policies and prevent goal misgener-
 346 alization in all the cases we studied, indicating there is still room for improvement between current
 347 methods and the theoretical ideal. These findings highlight MMER-based training as a promising
 348 approach to preventing goal misgeneralization.

349 A Proofs for theoretical results from Section 5

350 **Theorem 1** (MEV is susceptible to goal misgeneralization). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$,
351 a pair of normalized goals R, \tilde{R} , a proxy-distinguishing distribution shift $\langle \alpha, \beta, C, \Lambda^{\text{Train}}, \Lambda^{\text{Deploy}} \rangle$,
352 and an approximation threshold $\varepsilon \geq 0$. If $\varepsilon \geq \alpha$, then*

$$\exists \pi^{\text{MEV}} \in \Pi_{\varepsilon}^{\text{MEV}}(R, \Lambda^{\text{Train}}) \text{ such that } V^R(\pi^{\text{MEV}}; \Lambda^{\text{Deploy}}) < \max_{\pi \in \Pi} V^R(\pi; \Lambda^{\text{Deploy}}) - \beta C.$$

353 *Proof.* Construct policies $\pi^*, \tilde{\pi}^* \in \Pi$ that are optimal in all levels under R and \tilde{R} , respectively.
354 Assume $\varepsilon \geq \alpha$. We will put $\pi^{\text{MEV}} = \tilde{\pi}^*$. It remains to show (1) $\tilde{\pi}^* \in \Pi_{\varepsilon}^{\text{MEV}}(R, \Lambda^{\text{Train}})$ and (2)
355 $V^R(\tilde{\pi}^*; \Lambda^{\text{Deploy}}) < \max_{\pi \in \Pi} V^R(\pi; \Lambda^{\text{Deploy}}) - \beta C$. For (1), let $\Lambda_{\text{Distg.}}^{\text{Train}}, \Lambda_{\text{Ambig.}}^{\text{Train}} \in \Delta(\Theta)$ be Λ^{Train}
356 conditioned on the level being C -distinguishing or ambiguous, respectively. Then

$$\begin{aligned} V^R(\tilde{\pi}^*; \Lambda^{\text{Train}}) &= \alpha V^R(\tilde{\pi}^*; \Lambda_{\text{Distg.}}^{\text{Train}}) + (1 - \alpha) V^R(\tilde{\pi}^*; \Lambda_{\text{Ambig.}}^{\text{Train}}) && \text{(by Definition 3)} \\ &= \alpha V^R(\tilde{\pi}^*; \Lambda_{\text{Distg.}}^{\text{Train}}) + (1 - \alpha) V^R(\pi^*; \Lambda_{\text{Ambig.}}^{\text{Train}}) && \text{(by Definition 1)} \\ &\geq \alpha \cdot 0 + (1 - \alpha) V^R(\pi^*; \Lambda_{\text{Ambig.}}^{\text{Train}}) && \text{(since } V^R \geq 0) \\ &= V^R(\pi^*; \Lambda^{\text{Train}}) - \alpha V^R(\pi^*; \Lambda_{\text{Distg.}}^{\text{Train}}) && \text{(by Definition 3)} \\ &\geq V^R(\pi^*; \Lambda^{\text{Train}}) - \varepsilon \cdot 1. && \text{(since } \alpha \leq \varepsilon; V^R \leq 1) \end{aligned}$$

357 For (2), let $\Lambda_{\text{Distg.}}^{\text{Deploy}}, \Lambda_{\text{Ambig.}}^{\text{Deploy}} \in \Delta(\Theta)$ be defined similarly, from Λ^{Deploy} conditioned on the level being
358 C -distinguishing or ambiguous, respectively. Then

$$\begin{aligned} V^R(\tilde{\pi}^*; \Lambda^{\text{Deploy}}) &= \beta V^R(\tilde{\pi}^*; \Lambda_{\text{Distg.}}^{\text{Deploy}}) + (1 - \beta) V^R(\tilde{\pi}^*; \Lambda_{\text{Ambig.}}^{\text{Deploy}}) && \text{(by Definition 3)} \\ &= \beta V^R(\tilde{\pi}^*; \Lambda_{\text{Distg.}}^{\text{Deploy}}) + (1 - \beta) V^R(\pi^*; \Lambda_{\text{Ambig.}}^{\text{Deploy}}) && \text{(by Definition 1)} \\ &< \beta (V^R(\pi^*; \Lambda_{\text{Distg.}}^{\text{Deploy}}) - C) + (1 - \beta) V^R(\pi^*; \Lambda_{\text{Ambig.}}^{\text{Deploy}}) && \text{(by Definition 2)} \\ &= V^R(\pi^*; \Lambda^{\text{Deploy}}) - \beta C. && \text{(by Definition 3)} \quad \square \end{aligned}$$

359 **Theorem 2** (MMER is robust to goal misgeneralization). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a
360 pair of goals R, \tilde{R} , a proxy-distinguishing distribution shift $\langle \alpha, \beta, C, \Lambda^{\text{Train}}, \Lambda^{\text{Deploy}} \rangle$, and approxi-
361 mation thresholds $\varepsilon, \delta \geq 0$. Then*

$$\forall \pi^{\text{MMER}} \in \Pi_{\varepsilon, \delta}^{\text{MMER}}(R), \text{ we have } V^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) \geq \max_{\pi \in \Pi} V^R(\pi; \Lambda^{\text{Deploy}}) - \varepsilon - \delta.$$

362 *Proof.* Suppose $\pi^{\text{MMER}} \in \Pi_{\varepsilon, \delta}^{\text{MMER}}(R)$. Let $\Lambda^{\text{MMER}} \in \Delta(\Theta)$ such that $(\pi^{\text{MMER}}, \Lambda^{\text{MMER}})$ is an
363 (ε, δ) -equilibrium (Definition 5). Then we have the following bound on expected regret:

$$\begin{aligned} G^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) &\leq \max_{\Lambda \in \Delta(\Theta)} G^R(\pi^{\text{MMER}}; \Lambda) && (\Lambda^{\text{Deploy}} \in \Delta(\Theta)) \\ &\leq G^R(\pi^{\text{MMER}}; \Lambda^{\text{MMER}}) + \delta && (\Lambda^{\text{MMER}} \in \arg\text{-}\delta\text{-max}_{\Lambda \in \Delta(\Theta)} G^R(\pi^{\text{MMER}}; \Lambda)) \\ &\leq \min_{\pi \in \Pi} G^R(\pi; \Lambda^{\text{MMER}}) + \varepsilon + \delta. && (\pi^{\text{MMER}} \in \arg\text{-}\varepsilon\text{-min}_{\pi \in \Pi} G^R(\pi; \Lambda^{\text{MMER}})) \end{aligned}$$

364 We can convert this upper bound on expected regret to a lower bound on expected return:

$$\begin{aligned} V^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) &= \max_{\pi \in \Pi} V^R(\pi; \Lambda^{\text{Deploy}}) - G^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) && \text{(by equation 2)} \\ &\geq \max_{\pi \in \Pi} V^R(\pi; \Lambda^{\text{Deploy}}) - \min_{\pi \in \Pi} G^R(\pi; \Lambda^{\text{MMER}}) - \varepsilon - \delta. && \text{(by above bound)} \end{aligned}$$

365 The theorem follows, since, for any $\Lambda \in \Delta(\Theta)$, $\min_{\pi \in \Pi} G^R(\pi; \Lambda)$ vanishes by equation (2):

$$\begin{aligned} \min_{\pi \in \Pi} G^R(\pi; \Lambda) &= \min_{\pi \in \Pi} \left(\max_{\pi' \in \Pi} V^R(\pi'; \Lambda) - V^R(\pi; \Lambda) \right) && \text{(by equation 2)} \\ &= \max_{\pi' \in \Pi} V^R(\pi'; \Lambda) - \max_{\pi \in \Pi} V^R(\pi; \Lambda) && \text{(max term is constant wrt. } \pi) \\ &= 0. && \square \end{aligned}$$

Broader impact statement

Ngo et al. (2023) casts goal misgeneralization as a key risk mechanism for advanced deep learning systems, noting that techniques that improve capability robustness without preventing goal misgeneralization could *worsen* outcomes, since the system’s greater capabilities would then be devoted to the pursuit of an incorrect goal. Preventing this dangerous mode of generalization failure is a key challenge in assuring the safety of advanced RL agents.

In this section, we briefly note that minimax expected regret appears to be well-suited in principle to mitigating goal misgeneralization as deep learning systems become increasingly capable. This is because more generally capable deep learning systems should also be more capable regret-maximizing adversaries in particular. A more capable adversary will, in turn, be better at detecting or synthesizing rare, high-regret training situations, and then amplifying the training signal from these situations so as to induce correct generalization in an advanced deep RL agent.

Our work highlights training with the minimax expected regret (MMER) objective as a promising avenue for preventing goal misgeneralization. This objective has desirable theoretical properties, and we have found promising initial empirical results, though current MMER-based training techniques are not mature enough to prevent goal misgeneralization in all cases. However, as MMER-based training methods improve and as goal misgeneralization leads to more severe consequences, the ability of MMER-based training to mitigate goal misgeneralization should also improve.

Ultimately, we are hopeful that our work will instigate further research on the problem of goal misgeneralization, which remains a critical, open problem in the alignment and safe generalization of future advanced reinforcement learning agents.

References

- Matthew Barnett. A simple environment for showing mesa misalignment. Alignment Forum, September 2019. URL <https://www.alignmentforum.org/posts/AFdRGfYDWQqmkdhFq>.
- Michael Beukman, Samuel Coward, Michael Matthews, Mattie Fellows, Minqi Jiang, Michael Dennis, and Jakob Foerster. Refining minimax regret for unsupervised environment design. *arXiv preprint arXiv:2402.12284*, 2024.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/jax-ml/jax>.
- Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.
- Paul Christiano. Techniques for optimizing worst-case performance. AI Alignment (Blog), February 2018. URL <https://ai-alignment.com/techniques-for-optimizing-worst-case-performance-39eafec74b99>.
- Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1282–1289. PMLR, June 2019.
- Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 2048–2056. PMLR, 2020.

- 412 Michael Dennis, Natasha Jaques, Eugene Vinitisky, Alexandre Bayen, Stuart Russell, Andrew Critch,
413 and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised environment
414 design. In *Advances in Neural Information Processing Systems*, volume 33, pp. 13049–13061.
415 Curran Associates, Inc., 2020.
- 416 Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward, Yotam
417 Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu. IMPALA:
418 scalable distributed Deep-RL with importance weighted actor-learner architectures, 2018.
- 419 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel,
420 Matthias Bethge, and Felix A. Wichmann. Shortcut learning in deep neural networks. *Nature*
421 *Machine Intelligence*, 2:665–673, 2020. DOI: 10.1038/s42256-020-00257-z.
- 422 David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 423 Dylan Hadfield-Menell, Smitha Milli, Pieter Abbeel, Stuart Russell, and Anca Dragan. Inverse
424 reward design, 2017.
- 425 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
426 behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- 427 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
428 through world models. *arXiv preprint arXiv:2301.04104*, 2023.
- 429 Evan Hubinger. Towards an empirical investigation of inner alignment. Alignment Forum, Septem-
430 ber 2019. URL [https://www.alignmentforum.org/posts/2GycxikGnepJbxf](https://www.alignmentforum.org/posts/2GycxikGnepJbxfHT)
431 [HT](https://www.alignmentforum.org/posts/2GycxikGnepJbxfHT).
- 432 Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay, 2021.
- 433 Minqi Jiang, Michael Dennis, Jack Parker-Holder, Jakob Foerster, Edward Grefenstette, and Tim
434 Rocktäschel. Replay-guided adversarial environment design, 2022.
- 435 Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of zero-shot gener-
436 alisation in deep reinforcement learning. *Journal of Artificial Intelligence Research*, 76:201–264,
437 2023.
- 438 Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. Rma: Rapid motor adaptation for
439 legged robots. *arXiv preprint arXiv:2107.04034*, 2021.
- 440 Lauro Langosco Di Langosco, Jack Koch, Lee D Sharkey, Jacob Pfau, and David Krueger. Goal
441 misgeneralization in deep reinforcement learning. In *Proceedings of the 39th International Con-*
442 *ference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp.
443 12004–12019. PMLR, July 2022.
- 444 Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa,
445 Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training
446 group information. In *International Conference on Machine Learning*, pp. 6781–6792. PMLR,
447 2021.
- 448 Yecheng Jason Ma, William Liang, Hung-Ju Wang, Sam Wang, Yuke Zhu, Linxi Fan, Osbert Bas-
449 tani, and Dinesh Jayaraman. Dreureka: Language model guided sim-to-real transfer. *arXiv*
450 *preprint arXiv:2406.01967*, 2024.
- 451 Viktor Makoviychuk, Lukasz Wawrzyniak, Yunrong Guo, Michelle Lu, Kier Storey, Miles Macklin,
452 David Hoeller, Nikita Rudin, Arthur Allshire, Ankur Handa, et al. Isaac gym: High performance
453 gpu-based physics simulation for robot learning. *arXiv preprint arXiv:2108.10470*, 2021.
- 454 Fabio Muratore, Fabio Ramos, Greg Turk, Wenhao Yu, Michael Gienger, and Jan Peters. Robot
455 learning from randomized simulations: A review. *Frontiers in Robotics and AI*, 9:799893, 2022.

- Richard Ngo, Lawrence Chan, and Sören Mindermann. The alignment problem from a deep learning perspective, 2023.
- OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider, Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba, and Lei Zhang. Solving Rubik’s Cube with a robot hand. Preprint arXiv:1910.07113 [cs.LG], 2019.
- Jack Parker-Holder, Mingqi Jiang, Michael Dennis, Mikayel Samvelyan, Jakob Foerster, Edward Grefenstette, and Tim Rocktäschel. Evolving curricula with regret-based environment design, 2023.
- Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 3803–3810. IEEE, 2018.
- Alexander Rutherford, Michael Beukman, Timon Willi, Bruno Lacerda, Nick Hawes, and Jakob Foerster. No regrets: Investigating and improving regret approximations for curriculum discovery. In *Advances in Neural Information Processing Systems 37*, pp. 16071–16101. Curran Associates, 2024.
- Leonard J. Savage. The theory of statistical decision. *Journal of the American Statistical Association*, 46(253):55–67, 1951.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering Atari, Go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Rohin Shah, Vikrant Varma, Ramana Kumar, Mary Phuong, Victoria Krakovna, Jonathan Uesato, and Zac Kenton. Goal misgeneralization: Why correct specifications aren’t enough for correct goals, 2022.
- G. Starace. Addressing goal misgeneralization with natural language interfaces. Master’s thesis, Faculty of Science (FNWI), University of Amsterdam, 2023.
- Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017.
- Tu Trinh, Mohamad H Danesh, Nguyen X Khanh, and Benjamin Plaut. Getting by goal misgeneralization with a little help from a mentor. *arXiv preprint arXiv:2410.21052*, 2024.
- Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. *arXiv preprint arXiv:2408.14837*, 2024.
- Shengbo Wang, Nian Si, Jose Blanchet, and Zhengyuan Zhou. On the foundation of distributionally robust reinforcement learning. Preprint arxiv:2311.09018 [cs.LG], 2023.
- Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Ré. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations. *arXiv preprint arXiv:2203.01517*, 2022.
- Simon Zhuang and Dylan Hadfield-Menell. Consequences of misaligned AI. *Advances in Neural Information Processing Systems*, 33:15763–15773, 2020.

Supplementary Materials

The following content was not necessarily subject to peer review.

Contents

B	Expected regret identity for UMDPs	16
C	Partially observable environments	17
C.1	Generalizing the expected regret identity to partially observable environments . . .	17
C.2	Generalizing Definition 5 and Theorem 2 to partially observable environments . . .	18
D	Additional environment details	19
D.1	The CHEESE IN THE CORNER environment	19
D.2	The CHEESE ON A DISH environment	21
D.3	The KEYS AND CHESTS environment	23
E	Additional training details	27
E.1	Hyperparameters	27
E.2	Compute	27
F	Performance on ambiguous levels and with respect to the proxy goal	28
G	Visualizing performance on levels with different cheese positions	29
H	Experiments with different edit distributions	31
H.1	Elementary edits	31
H.2	ACCEL variants	31
H.3	Experimental results	32
I	The maximin expected value objective is susceptible to goal misgeneralization	34
I.1	Theoretical results	34
I.2	Training methods and experimental results	36
J	Experiments with increased training length	38
K	Experiments with a different distinguishing level generator	39
L	Experiments with different observations	40

530 B Expected regret identity for UMDPs

531 In this section, we prove equation (2) for UMDPs. Recall the following definitions from Section 3.

$$V^R(\pi; \theta) = \mathbb{E}_{s_0 \sim \mathcal{I}(\theta), a_t \sim \pi(\theta, s_t), s_{t+1} \sim \mathcal{T}(\theta, s_t, a_t)} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1}) \right] \quad (5)$$

$$V^R(\pi; \Lambda) = \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi; \theta)] \quad (6)$$

$$G^R(\pi; \theta) = \max_{\pi' \in \Pi} V^R(\pi'; \theta) - V^R(\pi; \theta) \quad (7)$$

$$G^R(\pi; \Lambda) = \mathbb{E}_{\theta \sim \Lambda} [G^R(\pi; \theta)] \quad (8)$$

532 In Section 3 we observe that, for UMDPs, we have the additional basic identity

$$G^R(\pi; \Lambda) = \max_{\pi' \in \Pi} V^R(\pi'; \Lambda) - V^R(\pi; \Lambda). \quad (\text{this is equation 2})$$

533 This is a nontrivial identity that does not hold for partially observable underspecified environments
 534 in which the level is not observable to the policy (see Appendix C.1). However, for policies that are
 535 conditioned on the level, the identity holds, as we now prove.

536 **Proposition 1** (Expected regret identity for UMDPs). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a goal*
 537 *R , and a level distribution $\Lambda \in \Delta(\Theta)$. Let Π be the set of all policies of the form $\pi : \Theta \times \mathcal{S} \rightarrow \Delta(\mathcal{A})$.*
 538 *Then we have*

$$G^R(\pi; \Lambda) = \max_{\pi' \in \Pi} V^R(\pi'; \Lambda) - V^R(\pi; \Lambda).$$

539 *Proof.* $G^R(\pi; \Lambda) = \mathbb{E}_{\theta \sim \Lambda} [G^R(\pi; \theta)]$ (by equation 8)

$$= \mathbb{E}_{\theta \sim \Lambda} \left[\max_{\pi' \in \Pi} V^R(\pi'; \theta) - V^R(\pi; \theta) \right] \quad (\text{by equation 7})$$

$$= \max_{\pi' \in \Pi} \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi'; \theta)] - \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi; \theta)] \quad (\text{by Proposition 2, below})$$

$$= \max_{\pi' \in \Pi} V^R(\pi'; \Lambda) - V^R(\pi; \Lambda). \quad (\text{by equation 6}) \quad \square$$

540 The above proof relies on Proposition 2, which says we can exchange expectation and maximization
 541 for the expected return since the policy is conditioned on the level.

542 **Proposition 2** (Expectation and maximization of expected return commute for UMDPs). *Consider*
 543 *an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a goal R , and a level distribution $\Lambda \in \Delta(\Theta)$. Let Π be the set of all*
 544 *policies of the form $\pi : \Theta \times \mathcal{S} \rightarrow \Delta(\mathcal{A})$. Then we have*

$$\mathbb{E}_{\theta \sim \Lambda} \left[\max_{\pi \in \Pi} V^R(\pi; \theta) \right] = \max_{\pi \in \Pi} \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi; \theta)].$$

545 *Proof.* (\geq): Note that this direction holds regardless of whether we condition policies on the level.

546 Let $\pi^* \in \arg \max_{\pi \in \Pi} \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi; \theta)]$. Then we have

$$\max_{\pi \in \Pi} \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi; \theta)] = \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi^*; \theta)] \leq \mathbb{E}_{\theta \sim \Lambda} \left[\max_{\pi \in \Pi} V^R(\pi; \theta) \right].$$

547 (\leq): Observe that, per equation (5), $V^R(\pi; \theta)$ depends only on π through $\pi(\theta, -) : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, that
 548 is, through the policy conditioned on the fixed level θ . Therefore we can construct a single policy that
 549 achieves the maximum expected return under all levels. For $\theta \in \Theta$, let $\pi_\theta^* \in \arg \max_{\pi \in \Pi} V^R(\pi; \theta)$.
 550 Then define $\pi^* : \Theta \times \mathcal{S} \rightarrow \Delta(\mathcal{A})$ such that for $\theta \in \Theta$, $s \in \mathcal{S}$, and $a \in \mathcal{A}$, $\pi^*(a | \theta, s) = \pi_\theta^*(a | \theta, s)$.
 551 By construction, we have $\pi_\theta^* \in \arg \max_{\pi \in \Pi} V^R(\pi; \theta)$ for all $\theta \in \Theta$. It follows that

$$\mathbb{E}_{\theta \sim \Lambda} \left[\max_{\pi \in \Pi} V^R(\pi; \theta) \right] = \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi^*; \theta)] \leq \max_{\pi \in \Pi} \mathbb{E}_{\theta \sim \Lambda} [V^R(\pi; \theta)]. \quad \square$$

552 C Partially observable environments

553 In this section, we generalize equation (2) and Theorem 2 to partially observable environments, with
 554 a slight modification of the bound to account for the fact that it may no longer be possible for any
 555 policy to achieve zero expected regret on a given distribution of levels.

556 Rather than defining underspecified partially observable MDPs in detail, we consider an arbitrary
 557 subset of the space of policies $\Phi \subseteq \Pi$. We can model partial observability by restricting to policies
 558 with tied outputs within any given partition of $\Theta \times \mathcal{S}$ into information sets.

559 We define the **expected restricted regret** of a policy $\pi \in \Phi$ in a level $\theta \in \Theta$ under a goal R based
 560 on the return of the best available policy within such a subset of policies:

$$G_{\Phi}^R(\pi; \theta) = \max_{\pi' \in \Phi} V^R(\pi'; \theta) - V^R(\pi; \theta). \quad (9)$$

561 As before, we lift this definition to a level distribution Λ by taking the expectation

$$G_{\Phi}^R(\pi; \Lambda) = \mathbb{E}_{\theta \sim \Lambda} [G_{\Phi}^R(\pi; \theta)]. \quad (10)$$

562 C.1 Generalizing the expected regret identity to partially observable environments

563 **Proposition 3** (Expected regret identity for partially observable environments). *Consider an UMDP*
 564 *$\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a goal R , a level distribution $\Lambda \in \Delta(\Theta)$, and a subset of policies $\Phi \subseteq \Pi$. Then*

$$G_{\Phi}^R(\pi; \Lambda) = \max_{\pi' \in \Phi} V^R(\pi'; \Lambda) - V^R(\pi; \Lambda) + \min_{\pi' \in \Phi} G_{\Phi}^R(\pi'; \Lambda).$$

565 *Proof.* Equivalently,

$$\begin{aligned} & G_{\Phi}^R(\pi; \Lambda) - \min_{\pi' \in \Phi} G_{\Phi}^R(\pi'; \Lambda) \\ &= \mathbb{E}_{\theta \sim \Lambda} [G_{\Phi}^R(\pi; \theta)] - \min_{\pi' \in \Phi} \mathbb{E}_{\theta \sim \Lambda} [G_{\Phi}^R(\pi'; \theta)] \quad (\text{by equation 10}) \\ &= \mathbb{E}_{\theta \sim \Lambda} \left[\max_{\pi' \in \Phi} V^R(\pi'; \theta) \right] - V^R(\pi; \Lambda) - \min_{\pi' \in \Phi} \left(\mathbb{E}_{\theta \sim \Lambda} \left[\max_{\pi'' \in \Phi} V^R(\pi''; \theta) \right] - V^R(\pi'; \Lambda) \right) \\ & \quad (\text{by equations 9 and 6}) \\ &= \mathbb{E}_{\theta \sim \Lambda} \left[\max_{\pi' \in \Phi} V^R(\pi'; \theta) \right] - V^R(\pi; \Lambda) - \mathbb{E}_{\theta \sim \Lambda} \left[\max_{\pi'' \in \Phi} V^R(\pi''; \theta) \right] + \max_{\pi' \in \Phi} V^R(\pi'; \Lambda) \\ &= \max_{\pi' \in \Phi} V^R(\pi'; \Lambda) - V^R(\pi; \Lambda). \quad \square \end{aligned}$$

566 Compared to equation (2) (Proposition 1), Proposition 3 includes the term $\min_{\pi \in \Phi} G_{\Phi}^R(\pi; \Lambda)$. This
 567 extra term represents the **irreducible (restricted expected) regret** for the level distribution Λ .
 568 Proposition 3 essentially says that the restricted expected regret for a level distribution can be de-
 569 composed into two components: (1) the shortfall in expected return compared to the optimal policy
 570 for the level distribution (as in equation 2, cf. the definition of regret for individual levels); plus
 571 (2) this irreducible regret.

572 Irreducible regret can arise when the level is partially observable to the policy. For example, consider
 573 a mixture of two levels with two disjoint sets of optimal policies. Suppose the level is not observed
 574 by the policy, so the policy has to choose actions without knowing whether it is in the first level
 575 or the second level. In each individual level, we define expected regret based on the performance
 576 of optimal policies for that level (these policies will naturally be the ones that *assume* they are in
 577 the appropriate level). However, since no single policy can perform optimally in both levels, the
 578 expected regret with respect to the mixture is always nonzero. This nonzero minimum expected
 579 regret is exactly the irreducible regret.

580 C.2 Generalizing Definition 5 and Theorem 2 to partially observable environments

581 We are now in position to generalize the results of Section 5.2 to partially observable environments.

582 **Definition 6** (Approximate MMER for partially observable environments). *Consider an UMDP*
 583 $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, *a goal* R , *approximation thresholds* $\varepsilon, \delta \geq 0$, *and a subset of policies*
 584 $\Phi \subseteq \Pi$. *Consider the two-player zero-sum game* $\langle \langle \Phi, \Delta(\Theta) \rangle, \langle -G_\Phi^R, G_\Phi^R \rangle \rangle$ (cf. Definition 5).
 585 *A pair* (π, Λ) *is a restricted* (ε, δ) -*equilibrium if* $\pi \in \arg\text{-}\varepsilon\text{-}\min_{\pi' \in \Phi} G_\Phi^R(\pi'; \Lambda)$ *and* $\Lambda \in$
 586 $\arg\text{-}\delta\text{-}\max_{\Lambda' \in \Delta(\Theta)} G_\Phi^R(\pi; \Lambda')$. *The restricted approximate equilibrium set is then*

$$\mathcal{E}_{\varepsilon, \delta}^{\text{MMER}}(R) = \{(\pi, \Lambda) \in \Phi \times \Delta(\Theta) \mid (\pi, \Lambda) \text{ is a restricted } (\varepsilon, \delta)\text{-equilibrium}\}.$$

587 **Theorem 3** (MMER is robust to goal misgeneralization in partially observable environments). *Con-*
 588 *sider an UMDP* $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, *a pair of goals* R, \tilde{R} , *a proxy-distinguishing distribution shift*
 589 $\langle \alpha, \beta, C, \Lambda^{\text{Train}}, \Lambda^{\text{Deploy}} \rangle$, *approximation thresholds* $\varepsilon, \delta \geq 0$, *and a subset of policies* $\Phi \subseteq \Pi$. *Then*
 590 *for all restricted* (ε, δ) -*equilibria* $(\pi^{\text{MMER}}, \Lambda^{\text{MMER}}) \in \mathcal{E}_{\varepsilon, \delta}^{\text{MMER}}(R)$, *we have*

$$V^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) \geq \max_{\pi \in \Phi} V^R(\pi; \Lambda^{\text{Deploy}}) - \varepsilon - \delta - g(\Lambda^{\text{MMER}}, \Lambda^{\text{Deploy}})$$

591 *where* $g(\Lambda^{\text{MMER}}, \Lambda^{\text{Deploy}}) = \min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{MMER}}) - \min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{Deploy}})$ *is the difference in*
 592 *irreducible expected restricted regret between* Λ^{MMER} *and* Λ^{Deploy} .

593 *Proof.* Let $(\pi^{\text{MMER}}, \Lambda^{\text{MMER}})$ be a restricted (ε, δ) -equilibrium. Following analogous steps to the
 594 proof of Theorem 2, we have an upper bound on the restricted expected regret by Definition 6,

$$G_\Phi^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) \leq \min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{MMER}}) + \varepsilon + \delta.$$

595 Once again, we can convert this upper bound on expected regret to a lower bound on expected return:

$$\begin{aligned} & V^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) \\ &= \max_{\pi \in \Phi} V^R(\pi; \Lambda^{\text{Deploy}}) - G_\Phi^R(\pi^{\text{MMER}}; \Lambda^{\text{Deploy}}) + \min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{Deploy}}) \quad (\text{by Proposition 3}) \\ &\geq \max_{\pi \in \Phi} V^R(\pi; \Lambda^{\text{Deploy}}) - \varepsilon - \delta - \min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{MMER}}) + \min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{Deploy}}) \quad (\text{by above bound}) \\ &= \max_{\pi \in \Phi} V^R(\pi; \Lambda^{\text{Deploy}}) - \varepsilon - \delta - g(\Lambda^{\text{MMER}}, \Lambda^{\text{Deploy}}). \quad \square \end{aligned}$$

596 Compared to Theorem 2, the performance bound in Theorem 3 has the additional term

$$g(\Lambda^{\text{MMER}}, \Lambda^{\text{Deploy}}) = \min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{MMER}}) - \min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{Deploy}}). \quad (11)$$

597 We make the following remarks:

- 598 1. The term $g(\Lambda^{\text{MMER}}, \Lambda^{\text{Deploy}})$ represents how much higher the irreducible regret of the MMER
 599 distribution is compared to the irreducible regret of the deployment distribution. This difference
 600 comes out of the performance guarantee because, in the MMER game, the agent has no incentive
 601 to improve regret on the deployment distribution once it is below the irreducible regret of the
 602 MMER distribution. This is a limitation of standard MMER that can be addressed by lexico-
 603 graphic refinement of the decision rule along the lines of Beukman et al. (2024).
- 604 2. Note $\min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{MMER}}) \in [G_\Phi^{\text{MMER}} - \varepsilon, G_\Phi^{\text{MMER}}]$ where $G_\Phi^{\text{MMER}} = G_\Phi^R(\pi^{\text{MMER}}; \Lambda^{\text{MMER}})$.
 605 We can estimate G_Φ^{MMER} at the end of MMER-based training by evaluating the agent's final policy
 606 in the agent's final level distribution. Thus we can estimate $\min_{\pi \in \Phi} G_\Phi^R(\pi; \Lambda^{\text{MMER}})$ in practice.
- 607 3. The bound holds for all restricted equilibria involving π^{MMER} . If one has access to multiple
 608 suitable Λ^{MMER} , one can take the best version of the performance bound.

609 D Additional environment details

610 In this appendix, we provide additional details about each environment, including details about pro-
 611 cedurally generating ambiguous and distinguishing levels, edit distributions, computation of maxi-
 612 mum level value for the oracle-latest estimator, and the origin of each environment.

613 D.1 The CHEESE IN THE CORNER environment

614 **Environment.** In this environment, levels are parameterized by a 13×13 wall layout, a mouse
 615 spawn position within this grid, and a cheese position within the grid. We require that the cheese
 616 position and the mouse spawn position are not equal, and moreover that they are not obstructed by
 617 walls. We do not require that there is an unobstructed path between them.

618 In the initial state, the mouse begins in the spawn position. The actions available to the agent are
 619 to attempt to move the mouse up, left, down, or right, which succeeds if the respective position is
 620 not obstructed by a wall or the edge of the grid. If the mouse position equals the cheese position,
 621 the mouse collects the cheese. The episode terminates when the cheese has been collected or after a
 622 maximum of 128 steps.

623 **Observations.** We represent states to the agent as a $15 \times 15 \times 3$ Boolean grid. The first of the
 624 three channels encodes the maze layout, including a border of width 1. The second channel one-hot
 625 encodes the position of the mouse. The third channel one-hot encodes the position of the cheese (if
 626 it has not been collected). All of the relevant information about the level and the state are encoded
 627 into this observation, therefore this environment is fully observable.

628 **True goal and proxy goal.** The training goal is for the mouse to collect the cheese. The reward
 629 function assigns +1 reward to transitions in which the mouse collects the cheese.

630 We also consider a proxy goal of navigating to the top left corner of the maze. This could be
 631 formulated as a reward function that assigns +1 reward the first time the mouse steps into the top
 632 left corner (this reward can be made Markovian by augmenting the state with a flag for whether the
 633 corner has previously been reached). Note that we never train with this proxy goal as the reward
 634 function.

635 **Level classification.** Given this environment and this pair of goals, we can classify levels according
 636 to the definitions in Section 4. Assume the discount factor is strictly between 0 and 1.

- 637 1. Consider a level for which the cheese is in the top-left corner of the maze, and this corner is reach-
 638 able from the mouse spawn position. Such a level is ambiguous: optimally pursuing the proxy
 639 goal implies following a shortest path to the corner, which is also a shortest path to collecting the
 640 cheese.
- 641 2. Consider a level for which the cheese is not in the top-left corner, but (1) the cheese is reachable
 642 from the mouse spawn position, (2) the top-left corner is reachable from the mouse spawn po-
 643 sition (other than via the cheese position), and (3) the top-left corner is not on any shortest path
 644 to the cheese from the mouse spawn position. Such a level is C -distinguishing for some positive
 645 C : Optimally pursuing the proxy goal must begin by taking a shortest path to the top-left corner,
 646 avoiding the cheese (to avoid early termination), which exists by (1). Thereafter, assuming the
 647 policy can condition on whether it has already visited the corner, any behavior is equally optimal
 648 under the proxy goal, and in particular such a policy may proceed to the cheese (which is reach-
 649 able from the top-left corner by (1) and (2)). However, by (3) the collection of the cheese will be
 650 delayed compared to taking a shortest path to the cheese position directly from the mouse spawn
 651 position. The extent of the delay and the discount factor determine C .

652 Other levels are not as straight-forward to classify. For example, if the top-left corner is on a shortest
 653 path to the cheese from the mouse spawn position, then some policies that are optimal for the proxy
 654 goal will be optimal for the true goal, but this will not be the case for all policies that are optimal
 655 for the proxy goal. Moreover, if the walls are configured such that the top-left corner is unreachable

from the mouse spawn position, then technically all policies are optimal according to the proxy goal (some or all of which may be optimal with respect to the true goal). Similarly, if the walls are configured such that the cheese is unreachable, then technically all policies are optimal with respect to the true goal (including all policies that are optimal with respect to the proxy goal). Note for determining reachability one should account for the episode termination conditions.

In practice, we normally use the simpler, approximately correct approach of checking whether the cheese is in the corner to determine whether a level is ambiguous or distinguishing. For example, this is the approach used in defining the proportion of distinguishing levels in the buffer. We note that all UED algorithms rapidly replace levels in which the cheese is unreachable with solvable levels in their buffer during the initial steps of training.

Procedural level generation. We construct two procedural level generators, $\Lambda_{\text{Ambig.}}, \Lambda_{\text{Distg.}} \in \Delta(\Theta)$, that are (approximately) concentrated on ambiguous and distinguishing levels, respectively.

- **Ambiguous level distribution** ($\Lambda_{\text{Ambig.}}$). We procedurally generate ambiguous levels as follows. For each position (except the top-left corner), we place a wall independently with probability 25%. We position the cheese in the top-left corner. We position the mouse spawn somewhere where there is not a wall, other than the top-left corner (we assume there is at least one such position).

All of the generated levels are ambiguous. However, we note that for many generated levels, the top-left corner is unreachable from the mouse spawn position. These levels are technically ambiguous (because all policies are optimal), and while this is not exactly the spirit of the definition, it is at least true that they do not provide a training signal in favor of the true goal over the proxy goal, which is enough for practical purposes.

- **Distinguishing level distribution** ($\Lambda_{\text{Distg.}}$). We procedurally generate *mostly* distinguishing levels as follows. For each position, we place a wall independently with probability 25%. We position the cheese somewhere where there is not a wall. We position the mouse spawn somewhere where there is not a wall, different from the cheese position.

The generated levels are mostly distinguishing levels (for varying ranges of C). It may arise that the cheese position is in the top-left corner, or moreover on a shortest path to the top-left corner from the mouse spawn position or vice versa, or that either the cheese position or the top-left corner or both are not reachable from the mouse spawn position. In such cases, the level may technically be either ambiguous or neither ambiguous nor distinguishing. Note that the effect of these pathological levels is only to reduce the training signal in favor of the true goal available during training (in the setting of a proxy-distinguishing distribution shift with low α).

Elementary edit distributions. The ACCEL training method additionally requires specifying an edit distribution used to sample “similar” levels for potential entry into the level buffer. In Appendix H, we discuss how we build our edit distributions from elementary edit distributions of three types: ambiguity preserving edits (used in ACCEL in the main text), biased ambiguity transforming edits, and unrestricted ambiguity transforming edits. For CHEESE IN THE CORNER, these elementary edit distributions are configured as follows.

1. **Ambiguity preserving edits.** Each ambiguity preserving edit either removes an existing wall, positions a new wall, or moves the mouse spawn position to a random unobstructed position (other than the cheese position). These edits don’t change the cheese position (they may change whether the cheese position is reachable from the mouse spawn position, or change whether the corner is on a shortest path to the cheese from the mouse spawn position).

2. **Biased ambiguity transforming edits.** Given a probability α , a biased ambiguity transforming edit randomizes the cheese position with probability α or places the cheese in the top-left corner with probability $1 - \alpha$. Note that when randomizing the cheese position, it’s possible that the cheese will be positioned in the top-left corner.

704 **3. Unrestricted ambiguity transforming edits.** An unrestricted ambiguity transforming edit ran-
 705 domizes the cheese position with probability 1. It’s possible that the cheese will be positioned in
 706 the top-left corner.

707 **Oracle maximum return.** Recall that the oracle-latest estimator (4) requires computing the max-
 708 imum expected return in a given level (under the true reward function). In this environment, an
 709 optimal policy follows any shortest path from the mouse spawn position to the cheese position in the
 710 graph representation of the maze layout. We compute the length of a shortest path using the Floyd–
 711 Warshall algorithm. Given a level $\theta \in \Theta$. Let R be the true reward function (+1 for collecting the
 712 cheese). If the shortest path has length $d \in \mathbb{N} \cup \{\infty\}$, then assuming a discount factor of $\gamma \in (0, 1)$,
 713 we have

$$\max_{\pi} V^R(\pi; \theta) = \gamma^{d-1}. \quad (12)$$

714 This includes the possibility that the cheese position is unreachable from the mouse spawn position,
 715 in which case $d = \infty$. Note that given mazes of this size, paths of length greater than the 128 are
 716 impossible, so we don’t have to consider the episode timeout.

717 **Origin.** CHEESE IN THE CORNER is an interpretation of the MAZE I environment from [Langosco](#)
 718 [et al. \(2022\)](#). [Langosco et al. \(2022\)](#) originally implemented the environment as a modification of the
 719 MAZE environment from OpenAI ProcGen ([Cobbe et al., 2020](#)), such that the cheese position could
 720 be restricted to a region surrounding the top-right corner (in Appendix K, we train on a distribution
 721 of levels with the cheese confined to top-left regions of varying size). [Hubinger \(2019\)](#) originated
 722 the idea of creating a navigation task with a location proxy as a potential means of inducing goal
 723 misgeneralization. Our environment uses an original implementation in JAX, and includes several
 724 changes from MAZE I including using a fixed size mazes and replacing the maze layout algorithm
 725 with a simpler algorithm based on random block placement.

726 D.2 The CHEESE ON A DISH environment

727 **Environment.** In this environment, levels are parameterized by a 13×13 wall layout, and positions
 728 within this grid for the mouse spawn, the dish, and the cheese. We require that the cheese position
 729 and the mouse spawn position are not equal, and that the dish position and the mouse spawn position
 730 are not equal, though the dish can be co-located with the cheese. Moreover, we require that none
 731 of the three positions are obstructed by walls. We do not require that there is an unobstructed path
 732 between the positions.

733 In the initial state, the mouse begins in the spawn position. The actions available to the agent are to
 734 attempt to move the mouse up, left, down, or right, which succeeds if the respective position is not
 735 obstructed by a wall or the edge of the grid. If the mouse position equals the cheese position, the
 736 mouse collects the cheese, and likewise for the dish. The episode terminates when the cheese *or* the
 737 dish has been collected or after a maximum of 128 steps.

738 **Observations.** We represent states to the agent as a $15 \times 15 \times (3 + D)$ Boolean grid where $D \in \mathbb{N}$.
 739 The first of the channels encodes the maze layout, including a border of width 1. The second
 740 channel one-hot encodes the position of the mouse. The third channel one-hot encodes the position
 741 of the cheese (if it has not been collected). The remaining D channels each one-hot encode the dish
 742 position (if the dish has not been collected). Redundantly coding the dish encourages the agent to
 743 learn a policy that is based on the dish position, eliciting a clearer case of goal misgeneralization.
 744 In our main experiments, $D = 6$. In Appendix L, we consider alternative values of D . All of the
 745 relevant information about the level and the state are encoded into this observation, therefore this
 746 environment is fully observable.

747 **True goal and proxy goal.** The training goal is for the mouse to collect the cheese. The reward
 748 function assigns +1 reward to transitions in which the mouse collects the cheese. We also consider
 749 a proxy goal of collecting the dish. This reward function assigns +1 reward to transitions in which
 750 the mouse collects the dish. Note that we never train with this proxy goal as the reward function.

Level classification. Given this environment and this pair of goals, we can classify levels according to the definitions in Section 4. Assume the discount factor is strictly between 0 and 1.

1. Consider a level for which the cheese and the dish are in the same position, and this position is reachable from the mouse spawn position. Such a level is ambiguous: optimally pursuing the proxy goal implies following the shortest path to collecting the dish, which is also the shortest path to collecting the cheese.
2. Consider a level for which (1) the cheese and the dish are not in the same position, (2) the cheese position and the dish position are both reachable from the mouse spawn position, including accounting for termination conditions. Such a level is C -distinguishing for some positive C : Optimally pursuing the proxy goal must involve taking the shortest path to the dish position that avoids the cheese (since collecting the cheese would terminate the episode before reaching the dish). This policy forfeits all available reward from the true goal, the amount of which depends on the length of the shortest path to the cheese position, avoiding the dish.

Other levels are not as straight-forward to classify. If the walls and the cheese are configured such that the dish position is not reachable from the mouse spawn position, then all policies are optimal according to the proxy goal (some or all of which may be optimal with respect to the true goal). Similarly, if the walls and the dish are configured such that the cheese position is not reachable from the mouse spawn position, then all policies are optimal according to the true goal (including all policies that are optimal with respect to the proxy goal).

In practice, we normally use the simpler, approximately correct approach of checking whether the cheese position equals the dish position to determine whether a level is ambiguous or distinguishing. For example, this is the approach used in defining the proportion of distinguishing levels in the buffer. We note that all UED algorithms rapidly replace levels in which the cheese is unreachable with solvable levels in their buffer during the initial steps of training.

Procedural level generation. We construct two procedural level generators, $\Lambda_{\text{Ambig.}}, \Lambda_{\text{Distg.}} \in \Delta(\Theta)$, that are (approximately) concentrated on ambiguous and distinguishing levels, respectively.

- **Ambiguous level distribution** ($\Lambda_{\text{Ambig.}}$). We procedurally generate ambiguous levels as follows. For each position, we place a wall independently with probability 25%. We position the mouse spawn somewhere where there is not a wall. We position the cheese and the dish somewhere where there is not a wall, other than the mouse spawn position (we assume there are at least two positions without walls).

All of the generated levels are ambiguous. However, we note that for some generated levels, the cheese position (which is also the dish position) is unreachable from the mouse spawn position. These levels are technically ambiguous (because all policies are optimal for both goals), and while this is not exactly the spirit of the definition, it is at least true that they do not provide a training signal in favor of the true goal over the proxy goal, which is enough for practical purposes.

- **Distinguishing level distribution** ($\Lambda_{\text{Distg.}}$). We procedurally generate *mostly* distinguishing levels as follows. For each position, we place a wall independently with probability 25%. We position the mouse spawn somewhere where there is not a wall. We position the cheese and the dish, independently, somewhere where there is not a wall, different from the mouse spawn position.

The generated levels are mostly distinguishing levels (for varying ranges of C). It may arise that the cheese position equals the dish position, or that either the cheese position or the dish position or both are not reachable from the mouse spawn position. In such cases, the level may technically be either ambiguous or neither ambiguous nor distinguishing. Note, like for CHEESE IN THE CORNER, that the effect of these pathological levels is only to reduce the training signal in favor of the true goal available during training (in the setting of a proxy-distinguishing distribution shift with low α).

Elementary edit distributions. The ACCEL training method additionally requires specifying an edit distribution used to sample “similar” levels for potential entry into the level buffer. In Ap-

pendix H, we discuss how we build our edit distributions from elementary edit distributions of three types: ambiguity preserving edits (used in ACCEL in the main text), biased ambiguity transforming edits, and unrestricted ambiguity transforming edits. For CHEESE ON A DISH, these elementary edit distributions are configured as follows.

1. **Ambiguity preserving edits.** Each ambiguity preserving edit either removes an existing wall, positions a new wall, or moves the mouse spawn position to a random unobstructed position (other than the cheese position or the dish position). These edits don't change the cheese position or the dish position. They may technically change whether these positions are reachable from the mouse spawn position.
2. **Biased ambiguity transforming edits.** Given a probability α , a biased ambiguity transforming edit randomizes the dish position, and then either independently randomizes the cheese position (with probability α) or positions the cheese at the new dish position (with probability $1 - \alpha$). Note that when independently randomizing the cheese position, it's possible that the cheese will be co-located with the dish.
3. **Unrestricted ambiguity transforming edits.** An unrestricted ambiguity transforming edit independently randomizes the cheese position and the dish position with probability 1. It's possible that the cheese and the dish will be co-located.

Oracle maximum return. Recall that the oracle-latest estimator (4) requires computing the maximum expected return in a given level (under the true reward function). In this environment, an optimal policy follows any shortest path from the mouse spawn position to the cheese position in the graph representation of the maze layout. If the dish position is not equal to the cheese position, the graph representation should treat the dish as an obstruction, since collecting the dish will terminate the episode before the cheese can be collected.

We compute the length of a shortest path in this graph using the Floyd–Warshall algorithm. Given a level $\theta \in \Theta$. Let R be the true reward function (+1 for collecting the cheese). If the shortest path has length $d \in \mathbb{N} \cup \{\infty\}$, then assuming a discount factor of $\gamma \in (0, 1)$, we have

$$\max_{\pi} V^R(\pi; \theta) = \gamma^{d-1}. \quad (13)$$

This includes the possibility that the cheese position is unreachable from the mouse spawn position without collecting the dish, in which case $d = \infty$. Note that given mazes of this size, paths of length greater than the 128 are impossible, so we don't have to consider the episode timeout.

Origin. This environment has a similar origin to that of CHEESE IN THE CORNER. It is an interpretation of the MAZE II environment from Langosco et al. (2022). Langosco et al. (2022) originally implemented the environment as a modification of the MAZE environment from OpenAI ProcGen (Cobbe et al., 2020), such that the cheese is replaced by a yellow diamond during training, and subsequently by a pair of objects (a red diamond and a yellow diagonal line) during evaluation. In their setup, the yellow diamond is rewarding because of its shape (diamond), rather than its color (yellow), so the intended extrapolation of the goal is for the policy to pursue the red diamond. However, Langosco et al. (2022) found that learned policies tend to pursue the yellow line instead. Since we are using Boolean grid observations rather than colored images, we redundantly code the dish to break the symmetry between the dish and the cheese, inducing a clearer case of goal misgeneralization.

D.3 The KEYS AND CHESTS environment

Environment. In this environment, levels are parameterized by a 13×13 wall layout, and positions within this grid for the mouse spawn, $k \leq 10$ keys, and $c \leq 10$ chests. We require that the mouse spawn position and the key and chest positions are distinct and are not obstructed by walls. We do not require that there are unobstructed paths between the positions.

In the initial state, the mouse begins in the spawn position. The actions available to the agent are to attempt to move the mouse up, left, down, or right, which succeeds if the respective position is

not obstructed by a wall or the edge of the grid. If the mouse position equals a key position, it collects the key, removing it from the maze and adding it to the mouse’s inventory. If the mouse position equals a chest position, assuming it has at least one key in its inventory, it collects the chest, removing it from the maze and removing the key from its inventory. The mouse can occupy the same position as a chest if it has an empty inventory. In a level with k keys and c chests, the episode terminates when the $\min(k, c)$ chests have been collected, or after a maximum of 128 steps.

Observations. We represent states to the agent as a $15 \times 15 \times 5$ Boolean grid. The first of the channels encodes the maze layout, including a border of width 1. The second channel one-hot encodes the position of the mouse. The third channel encodes the positions of all the as-yet-uncollected chests. The fourth channel encodes the positions of all as-yet-uncollected keys. The fifth channel encodes the mouse’s inventory, with one cell along the top row of the channel for each key. All of the relevant information about the level and the state are encoded into this observation, therefore this environment is fully observable.

True goal and proxy goal. The training goal is for the mouse to collect chests. The reward function assigns +1 reward to transitions in which the mouse collects a chest. Note that this reward function is not normalized like the ones we consider in the theory sections of the paper, but it is bounded and could easily be normalized.

Under this goal, collecting keys has no intrinsic value, but since collecting a chest requires collecting a key, keys have instrumental value. We also consider proxy goal that assign reward for collecting keys as well as chests. This goal could be modeled as a reward function that assigns, for example, $1 - \eta$ reward for collecting each key and η reward for collecting each chest, where $\eta \in (0, 1)$. Note that we never train with such a proxy goal as the reward function.

Level classification. Compared to the other environments, describing this environment in terms of the definitions in Section 4 is not as straight forward.

Optimal behavior in this environment involves collecting keys and chests in an order that depends on subtle tradeoffs driven by the exponential discounting. For the true goal, the optimal behavior is to collect as many chests as fast as possible, but it may make sense to make a brief detour to collect multiple keys if it slightly delays the collection of the first chest, as long as this sufficiently accelerates the collection of subsequent chests. The effect of switching to the proxy goal, which rewards key collection for its own sake, is to increase the incentive for the policy to take larger and larger detours to collect keys, and to incentivize collecting even more keys than there are reachable chests.

In our experiments, we mainly consider the following kinds of levels.

1. Consider a level in which there are $k \approx 3$ reachable keys, and $c \approx 10$ reachable chests. These levels are approximately ambiguous, in the sense that while the proxy goal incentivizes collecting keys earlier than optimal, it still incentivizes eventually collecting chests. For many key/chest layouts, pursuing the proxy goal entails similar behavior to pursuing the true goal.
2. Consider a level in which there are $k \approx 10$ reachable keys, and $c \approx 3$ reachable chests. These levels are mostly distinguishing, because, for many key/chest layouts, pursuing the proxy goal will usually involve a long detour to collect all reachable keys, delaying the collection of chests longer than is optimal.

These classifications depend on the exact key/chest layout. We note that keys and chests that are unreachable from the mouse spawn position have no influence on optimal behavior. However, as for the other environments, if there are no reachable keys or chests, this may affect the classification. In practice, we consider the total number of keys and chests in the level to classify it as ambiguous (3 keys and 10 chests) or distinguishing (10 keys and 3 chests), without accounting for reachability or the exact details of the key/chest layout.

Procedural level generation. We construct two procedural level generators, $\Lambda_{\text{Ambig.}}, \Lambda_{\text{Distg.}} \in \Delta(\Theta)$, that are (approximately) concentrated on ambiguous and distinguishing levels, respectively.

895 • **Ambiguous level distribution** ($\Lambda_{\text{Ambig.}}$). We procedurally generate *mostly* ambiguous levels as

896 follows. For each position, we place a wall independently with probability 25%. We then position
897 the mouse spawn, 3 keys, and 10 chests in distinct, unobstructed positions (assuming there are
898 enough positions). The generated levels are usually approximately ambiguous in the sense of the
899 classification system described above. It is possible that fewer than 3 keys and 10 chests will be
900 reachable from the mouse spawn position, and it is possible that the exact layout will lead to a
901 significant disadvantage for a policy that over-prioritizes key collection.

902 • **Distinguishing level distribution** ($\Lambda_{\text{Distg.}}$). We procedurally generate *mostly* distinguishing levels
903 in the same fashion, but with 10 keys and 3 chests instead. The generated levels are usually
904 distinguishing, along the lines of the classification system described above. It is possible that fewer
905 than 10 keys and 3 chests will be reachable from the mouse spawn position, and it is possible that
906 the exact layout will not sufficiently disincentivize key collection (for example, all keys positioned
907 on shortest paths to chests).

908 **Elementary edit distributions.** Once again, the ACCEL training method requires specifying an
909 edit distribution used to sample “similar” levels for potential entry into the level buffer. In Ap-
910 pendix H, we discuss how we build our edit distributions from elementary edit distributions of three
911 types: ambiguity preserving edits (used in ACCEL in the main text), biased ambiguity transforming
912 edits, and unrestricted ambiguity transforming edits. For KEYS AND CHESTS, these elementary edit
913 distributions are configured as follows.

914 1. **Ambiguity preserving edits.** Each ambiguity preserving edit either removes an existing wall,
915 positions a new wall, or moves the mouse spawn position or the position of one key or chest
916 to a random unobstructed, unoccupied position. These edits don’t change the number of keys
917 or chests. They may technically change whether these positions are reachable from the mouse
918 spawn position.

919 2. **Biased ambiguity transforming edits.** Given a probability α , a biased ambiguity transforming
920 edit sets the number of keys and chests to that of the distinguishing level generator (10 and 3)
921 with probability α , or that of the ambiguous level generator (3 and 10) with probability
922 $1 - \alpha$.

923 3. **Unrestricted ambiguity transforming edits.** An unrestricted ambiguity transforming edit is a
924 biased ambiguity transforming edit with α set to 50%.

925 **Oracle maximum return.** Recall, once more, that the oracle-latest estimator (4) requires comput-
926 ing the maximum expected return in a given level (under the true reward function). As described
927 above, an optimal policy in this environment collects keys and chests in the some order so as to
928 collect as many chests as fast as possible. In particular, noting that at most $m = \min(k, c)$ chests
929 can be collected in a given level (due to the requirement that a key must be expended to collect a
930 chest), and supposing that they are collected after steps $s_1, s_2, \dots, s_m \in \mathbb{N} \cup \{\infty\}$, the return is
931 given by

$$R(s_1, s_2, \dots, s_m) = \sum_{i=1}^m \gamma^{s_i-1} \quad (14)$$

932 where $\gamma \in (0, 1)$ is the discount factor.

933 Our approach to computing the optimal value is to enumerate a subset of paths through the network
934 of key/chest/mouse spawn positions that must contain an optimal path, to brute-force evaluate a
935 lower bound on the return of these paths, and to identify the greatest return lower bound as the
936 optimal return for the level. We explain the procedure in detail as follows.

937 1. We begin by enumerating a set of so-called *viable* collection sequences. Each collection se-
938 quences comprising an m -permutation of the k keys, an m -permutation of the c chests, and a
939 Dyck path of order m (a permutation of m keys and m chests such that each chest is preceded by
940 a corresponding key, cf. balanced parentheses). These three combinatorial objects jointly identify
941 a sequence in which particular keys and chests could be collected.

For example, suppose $k = 3$ and $c = 10$ and number the keys k_1, \dots, k_3 and the chests c_1, \dots, c_{10} . Suppose the key 3-permutation is $(3\ 1\ 2)$, the chest 3-permutation is $(1\ 6\ 4)$, and the Dyck path is $(k\ c\ k\ k\ c\ c)$. Then the corresponding collection sequence is $k_3, c_1, k_1, k_2, c_6, c_4$.

We note that the total number of viable collection sequences is

$$\underbrace{\binom{k}{m} \cdot m!}_{\text{key } m\text{-permutations}} \times \underbrace{\binom{c}{m} \cdot m!}_{\text{chest } m\text{-permutations}} \times \underbrace{\frac{1}{m+1} \binom{2m}{m}}_{\text{order } m \text{ Dyck paths}}. \quad (15)$$

Since, in our setup, either $k = 3$ and $c = 10$ or vice versa, $m = 3$ and (15) evaluates to 21,600.

2. We evaluate a lower bound on the return of each viable collection sequence as follows. First we compute all-pairs shortest path distances for the mouse spawn position, each key position, and each chest position, using the Floyd–Warshall algorithm. We represent unreachable pairs with a distance of ∞ . We then simulate each sequence, computing the step counts required for each collection as the cumulative sum of pairwise shortest path distances for each transition along the sequences, starting at the mouse spawn position. We round any distances above 128 up to ∞ . We then use the step counts of each of the m chest collections in the expression (14).

For example, consider the collection sequence described earlier, $k_3, c_1, k_1, k_2, c_6, c_4$. Let $D(p, q)$ represent the shortest path distance between the positions of objects p and q . Then we set $s_{c_1} = D(\text{spawn}, c_1)$, $s_{k_3} = s_{c_1} + D(c_1, k_3)$, and so on until $s_{c_4} = s_{c_6} + D(c_6, c_4)$. (If any of these step counts pass the timeout of 128, we round them up to ∞ to account for termination.) Finally, we compute the lower bound on the return of this sequence as $\gamma^{s_{c_1}-1} + \gamma^{s_{c_3}-1} + \gamma^{s_{c_4}-1}$.

We note that if a viable collection sequence ever involves collecting a key or chest that is unreachable from the mouse spawn position, then the step count for this collection and all subsequent collections will be infinite and thus this neither this collection nor subsequent collections will contribute to the return lower bound.

3. We take the greatest return lower bound across all viable collection sequences as the maximum return achievable in the level. We argue that this equals the maximum return as follows. Following a shortest path between a given pair of positions may involve crossing over other keys and chests. This can have any of several effects that invalidate the collection sequence, including (1) collecting keys or chests that appear later in the sequence earlier than planned, (2) expending keys before they are intended to be spent to collect a chest later in the sequence, or (3) terminating the episode early due to collecting the maximum available number of chests before finishing the sequence. However, such disruptions only ever *increase* the return. Moreover, for each disrupted sequence, there is a viable collection sequence that represents the actual order of keys and chests, except for accounting for the case where more than 3 keys are collected (which has no affect on the return as long as they are on the shortest paths to the necessary chests). It follows that for the viable collection sequence with the highest return, there are no such disruptions, and the return lower bound is tight.

We note that the set of viable collection sequences could be further filtered by eliminating (or never enumerating) sequences involving unreachable keys or chests. However, we use the above approach for simplicity and uniformity. In particular, since the above approach yields a fixed computational graph given values of k, c , we can enumerate the 21,600 viable sequences once at compile time and accelerate the brute-force evaluation step for a particular level (or batch of levels) using JAX. To handle a mixture of levels with $(k, c) = (3, 10)$ and levels with $(k, c) = (10, 3)$, we simply evaluate both ways and then dynamically keep the appropriate result for each level.

Origin. This environment is an interpretation of the KEYS AND CHESTS environment from [Langosco et al. \(2022\)](#). [Langosco et al. \(2022\)](#) originally implemented the environment as a modification of the HEIST environment from OpenAI ProcGen. [Barnett \(2019\)](#) originated the idea of creating a distribution shift from a navigation environment in which keys are scarce to one in which keys are not scarce. We implemented this environment in JAX, making several changes similar to those for the other environments.

E Additional training details

E.1 Hyperparameters

Table E.1: Hyperparameters used for all methods and environments.

Parameter	Value	Notes/exceptions
<i>Rollouts</i>		
# parallel environments	256	
Rollout length	128	
# environment steps per cycle	32,768	(# parallel environments \times rollout length)
Discount factor, γ	0.999	
<i>GAE</i>		
λ_{GAE}	0.95	
<i>PPO loss function</i>		
Clip range	0.1	
Value clipping	yes	
Critic coefficient	0.5	
Entropy coefficient	1e-3	Or, 1e-2 for KEYS AND CHESTS
<i>PPO updates</i>		
Epochs per cycle	5	
Minibatches per epoch	4	
Max gradient norm	0.5	
Adam learning rate	5e-5	
Learning rate schedule	constant	
<i>UED configuration</i>		N/A for DR
Replay rate		
PLR [⊥]	0.33	On average, 1 replay cycle per 2 generate cycles
ACCEL	0.5	On average, 1 replay cycle per 1 generate cycle (1 edit cycle immediately follows each replay cycle)
Buffer size	4096	
Prioritization method	rank	
Temperature	0.1	
Staleness coefficient	0.1	
# elementary edits per level, n	12	N/A for PLR [⊥]

E.2 Compute

We perform each training run on a single NVIDIA A100 80GB GPU. For CHEESE IN THE CORNER and CHEESE ON A DISH, each training run takes around 40 minutes (DR) or 80 minutes (UED methods). For KEYS AND CHESTS, each training run takes around 60 minutes (DR) or 110 minutes (UED methods). UED methods take longer than DR because UED methods require sampling additional rollouts for refining the buffer (the number of environment steps used for PPO updates is held constant across all methods). KEYS AND CHESTS runs are longer than others because we trained each method for 400 million steps instead of 200 million steps in this environment.

The experiments discussed in Section 7 took a total of 1.2k GPU hours. The additional experiments discussed in Appendices H, I, J, K, and L took, respectively, totals of 1.2k GPU hours; 500 GPU hours; 400 GPU hours (not counting the first 200 million environment steps used for training, which we counted with Section 7); 210 GPU hours; and 350 GPU hours.

1003 F Performance on ambiguous levels and with respect to the proxy goal

1004 In Section 7, we investigate which training distributions and methods led to good performance on
 1005 distinguishing levels. We claim that when the performance is low, this is an instance of goal misgen-
 1006 eralization, and therefore when the performance is high, goal misgeneralization has been prevented.

1007 In order to justify this claim, we also check that the poor performance is explained primarily by
 1008 the policy pursuing a proxy goal on distinguishing environments, rather than the policy behaving
 1009 incapably in ambiguous or distinguishing environments. Figure F.1 shows (1) return on ambiguous
 1010 levels and (2) proxy return on distinguishing levels for each training configuration.

1011 Figure F.1(top row) shows that the learned policies perform capably in ambiguous levels. For
 1012 CHEESE IN THE CORNER and CHEESE ON A DISH, all training methods achieve high return on
 1013 ambiguous levels for all training distributions. For KEYS AND CHESTS, this is the case for all train-
 1014 ing distributions except in $\alpha = 1$ baseline. Note that in this edge case, ambiguous levels (with few
 1015 keys and many chests) are never seen in training (cf. the case $\alpha = 0$ on distinguishing levels).

1016 Figure F.1(bottom row) shows the proxy return achieved by learned policies on distinguishing levels.
 1017 Note that we never use the proxy goal for training. Here we simply evaluate the policies according
 1018 to each environment’s respective proxy reward function. In particular, for CHEESE IN THE CORNER,
 1019 we use a reward function that assigns +1 reward the first time the mouse reaches the corner. For
 1020 CHEESE ON A DISH, we use a reward function that assigns +1 reward if the mouse collects the dish.
 1021 For KEYS AND CHESTS, it’s difficult to define a proxy goal (see Appendix D.3), here we report the
 1022 average number of keys in the mouse’s inventory throughout the rollouts (note that distinguishing
 1023 levels have at most 3 reachable chests, and so carrying more than three keys suggests the agent is
 1024 over-prioritizing key collection). The trends in proxy return mirror the trends in true return displayed
 1025 in Figure 3, suggesting that our learned policies retain enough capabilities in distinguishing levels
 1026 to pursue the proxy goal—a case of goal misgeneralization.

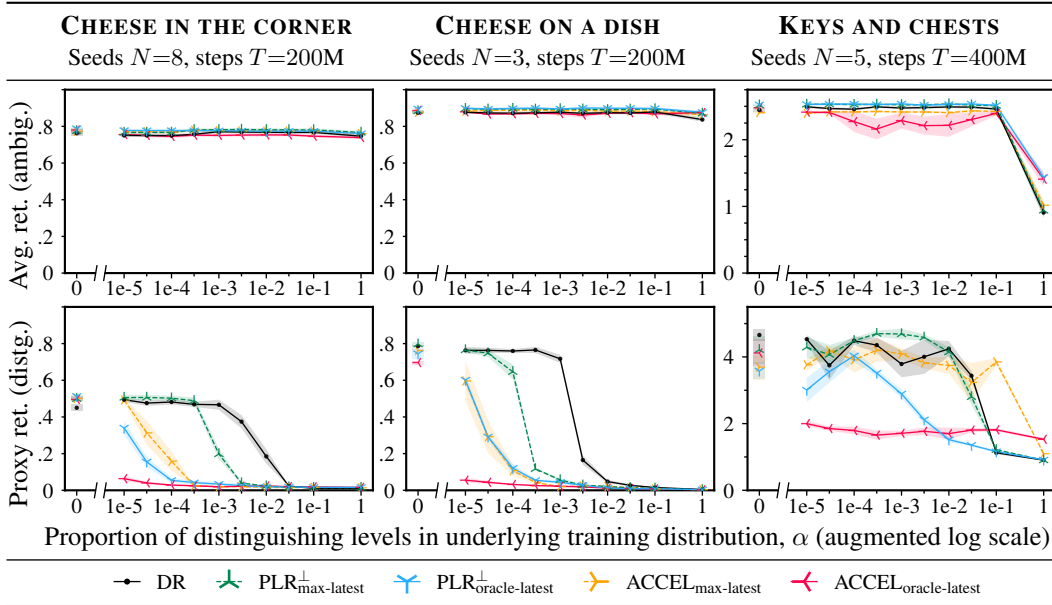


Figure F.1: **Performance on ambiguous levels and with respect to the proxy goal.** Each policy is trained on T environment steps using the indicated training method with underlying training distribution $\Lambda_{\alpha}^{\text{Train}}$. (*1st row*): Average return over 512 steps for an evaluation batch of 256 ambiguous levels. (*2nd row*): Average *proxy* return over 512 steps for an evaluation batch of 256 distinguishing levels. Note that the proxy goal is never used for training. (*Both*): Mean over N seeds, shaded region is one standard error. Note the split axes used to show zero on the log scale.

1027 G Visualizing performance on levels with different cheese positions

1028 For each training configuration (training distribution, training method) studied in the main text for
 1029 CHEESE IN THE CORNER, we save the policy from the end of training for the first of 8 training seeds.
 1030 In order to visualize the robustness of these policies to varying changes in the cheese position, we
 1031 create a batch of 122 levels with a shared wall layout and a fixed mouse position, but where each level
 1032 in the batch has a different cheese position. For each level, we sample 512 environment steps from
 1033 each policy and compute the average per-episode return as a measure of the policy’s performance in
 1034 that level. This gives us a vector of 122 average return values for each policy (one for each level),
 1035 which we visualize in a policy-specific heatmap such that the average return of the policy in a level
 1036 with the cheese in position (i, j) is indicated by the color of the cell (i, j) in the heatmap. For context
 1037 we overlay the wall layout and the mouse spawn position (note that we do not consider levels with
 1038 cheese positions that would coincide with a wall or with the mouse spawn position).

1039 Heatmaps for each method and training distribution follow in Figures G.1 and G.2. We see a rough
 1040 progression whereby for more advanced algorithms or higher α , the agent is robust to a greater
 1041 proportion of cheese positions. There are also instances of “blind spots” indicating cheese positions
 1042 for which certain methods are not robust, indicating that these training methods do not produce
 1043 perfectly robust policies.

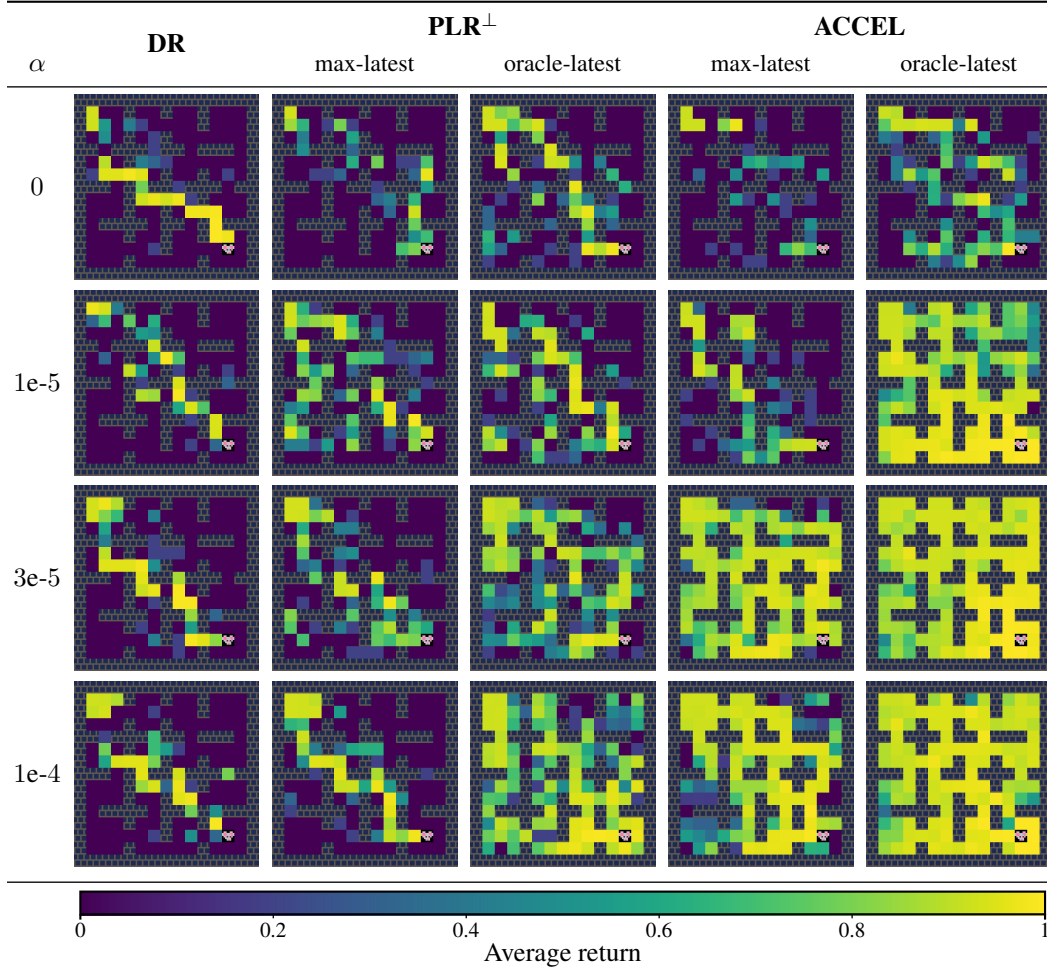


Figure G.1: Heatmap visualizations (part 1 of 2). See Figure 5 and Appendix G for details.

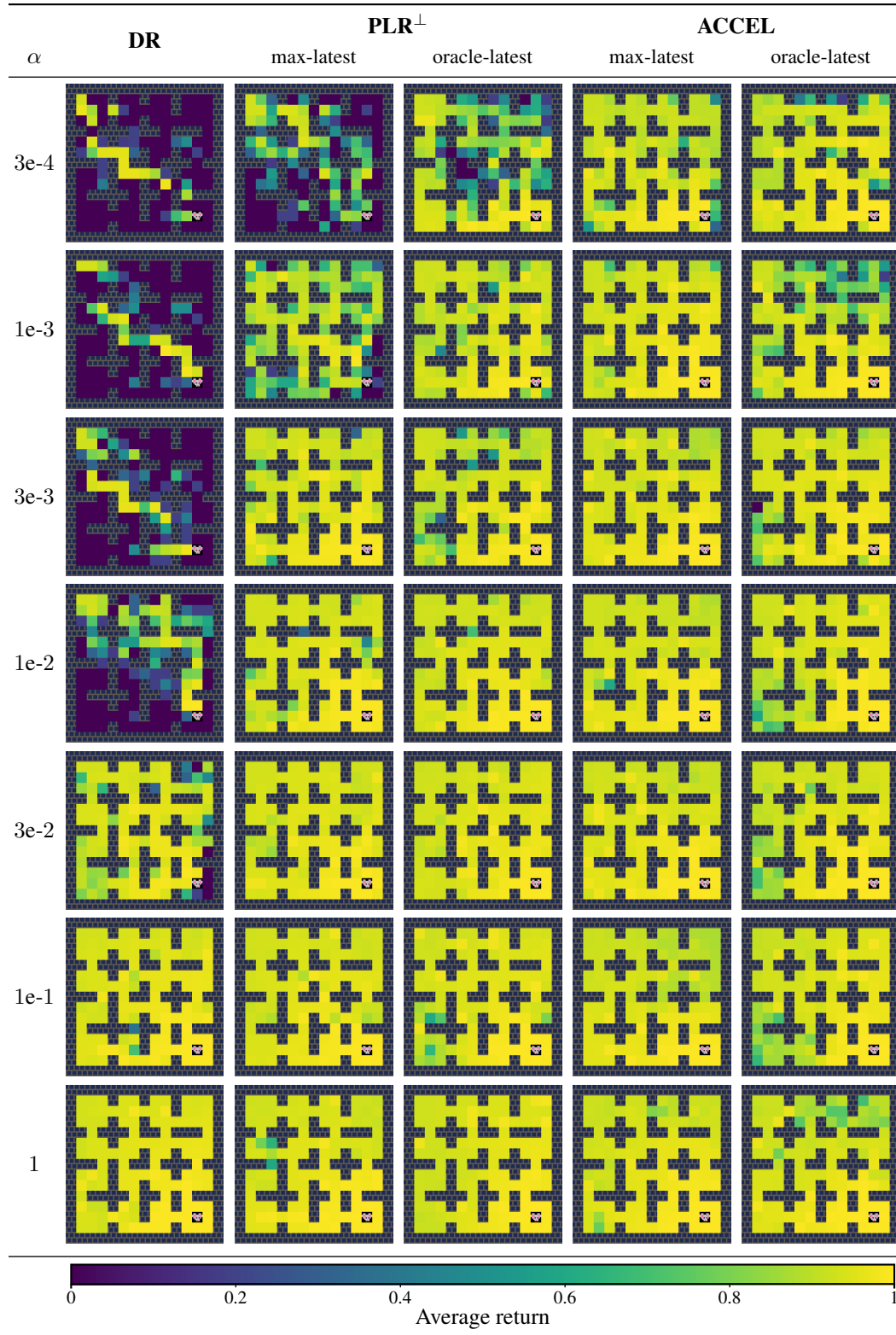


Figure G.2: Heatmap visualizations (part 2 of 2). See Figure 5 and Appendix G for details.

1044 H Experiments with different edit distributions

1045 As discussed in Section 6.2, ACCEL requires additionally specifying an **edit distribution** for mu-
 1046 tating levels in the buffer. In this appendix, we explore the effect of different edit distributions on the
 1047 ability for ACCEL to mitigate goal misgeneralization.

1048 H.1 Elementary edits

1049 We assemble our edit distributions by sampling a sequence of **elementary edits** of the following
 1050 three kinds.

- 1051 • **Ambiguity preserving edits.** These edits change the level without changing whether the level
 1052 is ambiguous or distinguishing. For example, in CHEESE IN THE CORNER, such an edit may
 1053 randomly toggle a wall or move the mouse’s starting position, but would not change the location
 1054 of the cheese.
- 1055 • **Biased ambiguity transforming edits.** These edits transform the level into a distinguishing level
 1056 with probability α or an ambiguous level with probability $1 - \alpha$, where α is the proportion of distin-
 1057 guishing levels in the underlying training distribution. For example, in CHEESE IN THE CORNER,
 1058 a biased ambiguity transforming edit may randomize the cheese position with probability α or
 1059 move it to the corner with probability $1 - \alpha$.
- 1060 • **Unrestricted ambiguity transforming edits.** These edits transform the position of the cheese
 1061 or the number of keys and chests uniformly at random given the level parameterization. For
 1062 CHEESE IN THE CORNER and CHEESE ON A DISH, the cheese moves to a random position in the
 1063 maze, which usually results in a distinguishing level. For KEYS AND CHESTS, we flip a coin to
 1064 make either keys or chests sparse, and the other type of object dense.

1065 We document the elementary edit distributions for each environment in full detail in Appendix D.

1066 H.2 ACCEL variants

1067 In these terms, we list the ACCEL variant considered in the main text along with three additional
 1068 variants of ACCEL with different edit distributions. We fix a hyperparameter n , the number of
 1069 elementary edits to apply to each level (we use $n = 12$).

- 1070 1. **Identity ACCEL** ($\text{ACCEL}^{\text{identity}}$, simply “ACCEL” in main text). We make a sequence of n
 1071 random edits, all of which are ambiguity preserving, resulting in the sequence of edits itself
 1072 being ambiguity preserving.
- 1073 2. **Constant ACCEL** ($\text{ACCEL}^{\text{constant}}$). We make $n - 1$ random ambiguity preserving edits, followed
 1074 by one biased ambiguity transforming edit. Applying this operation to any distribution of levels
 1075 results in a distribution with the same proportion of ambiguous and distinguishing levels as the
 1076 underlying training distribution.
- 1077 3. **Binomial ACCEL** ($\text{ACCEL}^{\text{binomial}}$). We make a sequence of n random edits, each independently
 1078 chosen to be either ambiguity preserving (with probability $1 - 1/n$) or else biased ambiguity
 1079 transforming. If the sequence has only ambiguity preserving edits (probability $(1 - 1/n)^n$) then
 1080 it is ambiguity preserving (like $\text{ACCEL}^{\text{identity}}$), otherwise the output is ambiguous with the same
 1081 probability as a level sampled from the underlying training distribution (like $\text{ACCEL}^{\text{constant}}$).
- 1082 4. **Unrestricted ACCEL** ($\text{ACCEL}^{\text{unrestricted}}$). We make a sequence of $n - 1$ random ambiguity pre-
 1083 serving edits, followed by one unrestricted ambiguity transforming edit. Applying this operation
 1084 to any distribution of levels results in a distribution with a proportion of distinguishing levels that
 1085 is independent of the parameter α that restricts access to distinguishing levels in the underlying
 1086 training distribution.

1087 $\text{ACCEL}^{\text{constant}}$ simulates restricted access to distinguishing levels. This variant is able to introduce
 1088 new distinguishing levels through edit operations, however, its ability to replicate existing distin-

guishing levels is limited, since every time it edits a level, the mutated level reverts to an ambiguous level with probability $1 - \alpha$.

ACCEL^{identity} (the variant studied in the main text) simulates a different kind of restriction, whereby we don't allow edits to turn ambiguous levels into distinguishing levels. However, we do allow edits to create new distinguishing levels by making similar copies of existing distinguishing levels in the buffer. Through this mechanism, ACCEL^{identity} can rapidly populate the buffer with distinguishing levels, providing this variant with additional capacity to amplify the small training signal from distinguishing levels (beyond simply curating levels, like in PLR[⊥] or ACCEL^{constant}).

ACCEL^{binomial} samples elementary edits in a different way. The number of biased ambiguity transforming edits included in the sequence follows a binomial distribution. This means that with around 35% probability, no biased ambiguity transforming edits will be applied, and the edit will resemble an edit from ACCEL^{identity}. Otherwise, with around 65% probability, at least one biased ambiguity transforming edit will be applied, and the overall edit will resemble one from ACCEL^{constant}. We thus expect the performance of this variant to be somewhere between that of ACCEL^{constant} and ACCEL^{identity}.

ACCEL^{unrestricted} is a baseline that simulates a situation where the edit distribution can be designed to explore the space of levels completely independently of the training distribution. We expect this variant to be able to obtain much stronger performance comparable to using $\alpha = 1$ in the training distribution, even when training with $\alpha = 0$.

H.3 Experimental results

We train with the three new variants and compare performance to DR and ACCEL^{identity} (from the main text). We report the results in Figure H.1 (oracle-latest regret estimator) and Figure H.2 (max-latest regret estimator). Note that we did not run ACCEL^{constant} with the oracle-latest regret estimator for CHEESE ON A DISH, or ACCEL^{constant} with the max-latest regret estimator for KEYS AND CHESTS.

Our results are in line with our central claim, that more advanced UED methods are more capable of mitigating goal misgeneralization.

- As predicted, ACCEL^{constant} achieves lower performance than ACCEL^{identity}. This could be explained by the greater flexibility with which ACCEL^{identity} can amplify distinguishing levels through edits.

- Moreover, ACCEL^{binomial} achieves performance somewhere between that of ACCEL^{constant} and ACCEL^{identity}.

- As predicted, ACCEL^{unrestricted} is able to populate the buffer with distinguishing levels regardless of the training distribution, even when $\alpha = 0$, both with the oracle-latest regret estimator and with the max-latest regret estimator.

ACCEL^{binomial} with max-latest estimation achieves the same low performance as ACCEL^{identity} in KEYS AND CHESTS (as observed for ACCEL^{identity} in Section 7).

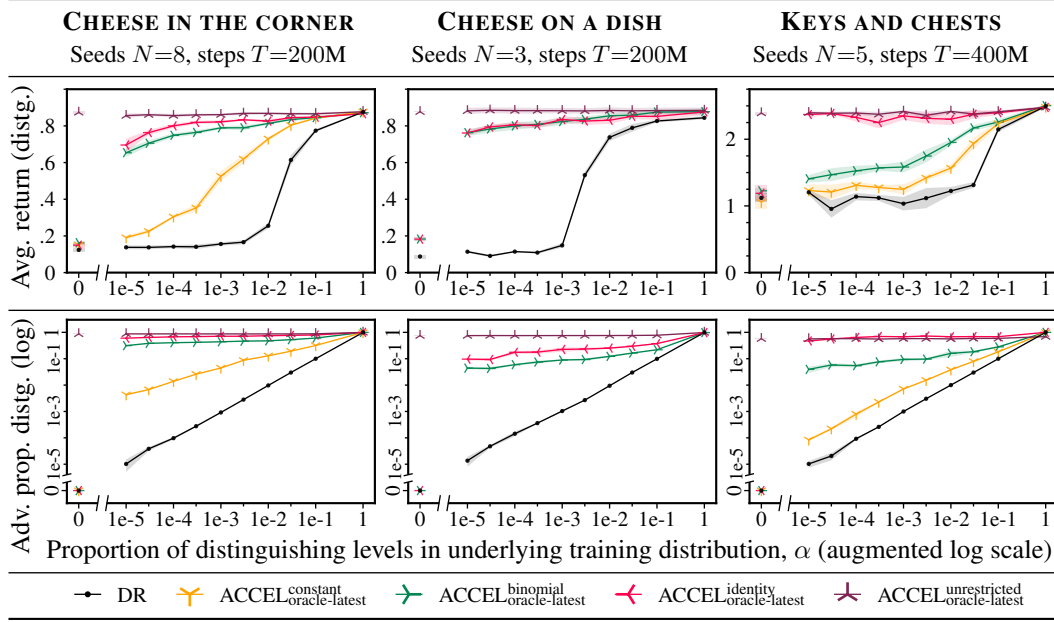


Figure H.1: ACCEL variants with oracle-latest estimator. Each policy is trained on T environment steps, using the indicated algorithm, with underlying training distribution $\Lambda_{\alpha}^{\text{Train}}$. (1st row): Average return over 512 steps for an evaluation batch of 256 distinguishing levels (cf. Figure 3). (2nd row): The proportion of distinguishing levels sampled from the adversary’s buffer across training (cf. Figure 4). (Both): Mean over N seeds, shaded region is one standard error. Note the split axes used to show zero on the log scale.

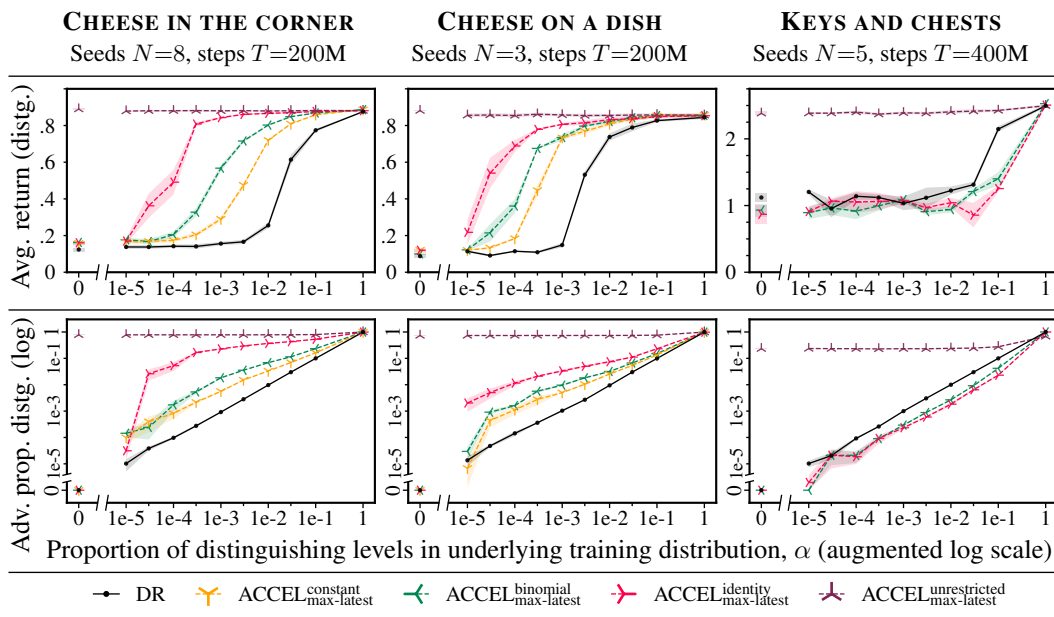


Figure H.2: ACCEL variants with max-latest estimator. Each policy is trained on T environment steps, using the indicated algorithm, with underlying training distribution $\Lambda_{\alpha}^{\text{Train}}$. (1st row): Average return over 512 steps for an evaluation batch of 256 distinguishing levels (cf. Figure 3). (2nd row): The proportion of distinguishing levels sampled from the adversary’s buffer across training (cf. Figure 4). (Both): Mean over N seeds, shaded region is one standard error. Note the split axes used to show zero on the log scale.

1126 I The maximin expected value objective is susceptible to goal 1127 misgeneralization

1128 In this section, we discuss an alternative strategy to minimax expected regret for selecting robust
1129 policies. Namely, we consider the **maximin expected value** (MMEV) training objective, whereby
1130 one seeks a policy that achieves the highest possible expected value (return) on a worst-case level
1131 distribution $\Lambda^{\text{MMEV}} \in \arg \min_{\Lambda' \in \Delta(\Theta)} V^R(\pi; \Lambda')$.

1132 [Dennis et al. \(2020\)](#) argues that the MMEV objective fails to induce robustness in an environment
1133 where the optimal value of each level varies within the level space. This is because the agent has
1134 no incentive to improve performance in any level above the maximum performance in worst-case
1135 levels. This same obstacle prevents MMEV from inducing robustness to goal misgeneralization,
1136 even though a policy that pursues a proxy goal is distinguishing levels will achieve low return in
1137 these levels.

1138 In this appendix, we show theoretically that MMEV-based methods allow for goal misgeneralization
1139 in environments with levels with low maximum value. Moreover, we show empirically that MMEV-
1140 based training methods fail to mitigate goal misgeneralization in our environments. Indeed, they fail
1141 to produce policies that perform capably in any levels.

1142 I.1 Theoretical results

1143 **Definition 7** (Approximate MMEV). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$, a goal R , and approx-
1144 imation thresholds $\varepsilon, \delta \geq 0$. Consider the two-player zero-sum game $\langle \langle \Pi, \Delta(\Theta) \rangle, \langle V^R, -V^R \rangle \rangle$,
1145 where an agent plays a policy $\pi \in \Pi$ and an adversary plays a level distribution $\Lambda \in \Delta(\Theta)$,
1146 aiming to maximize or minimize $V^R(\pi; \Lambda)$ respectively. A pair (π, Λ) is an (ε, δ) -equilibrium if
1147 $\pi \in \arg\text{-}\varepsilon\text{-}\max_{\pi' \in \Pi} V^R(\pi'; \Lambda)$ and $\Lambda \in \arg\text{-}\delta\text{-}\min_{\Lambda' \in \Delta(\Theta)} V^R(\pi'; \Lambda')$. The approximate MMEV
1148 policy set is then*

$$\Pi_{\varepsilon, \delta}^{\text{MMEV}}(R) = \{ \pi \in \Pi \mid \exists \Lambda \in \Delta(\Theta) \text{ such that } (\pi, \Lambda) \text{ is an } (\varepsilon, \delta)\text{-equilibrium} \},$$

1149 and the approximate MMEV level distribution set is

$$\Delta(\Theta)_{\varepsilon, \delta}^{\text{MMEV}}(R) = \{ \Lambda \in \Delta(\Theta) \mid \exists \pi \in \Pi \text{ such that } (\pi, \Lambda) \text{ is an } (\varepsilon, \delta)\text{-equilibrium} \}.$$

1150 Crucially, the MMEV adversary only cares to minimize the return of the agent. Given this, the
1151 equilibrium achieved is the one where the adversary plays levels in which *any* agent would achieve
1152 the minimum return possible.

1153 Let's first define some additional machinery we will need:

1154 **Definition 8** (α -minimum achievable return). *Given a level $\theta \in \Theta$ and a threshold α , define*

$$c(\theta, \alpha) = \min V_{\alpha}^R(\theta)$$

1155 where $V_{\alpha}^R(\theta) = \{ x \in \mathbb{R} \mid \exists \pi \in \Pi, \text{ such that } V^R(\pi; \theta) \geq \alpha \text{ and } V^R(\pi; \theta) = x \}$.

1156 **Theorem 4** (MMEV is susceptible to goal misgeneralization). *Consider an UMDP $\langle \Theta, \mathcal{A}, \mathcal{S}, \mathcal{I}, \mathcal{T} \rangle$,
1157 a pair of normalized goals R, \tilde{R} , a proxy-distinguishing distribution shift $\langle \alpha, \beta, C, \Lambda^{\text{Train}}, \Lambda^{\text{Deploy}} \rangle$,
1158 and approximation thresholds $\varepsilon, \delta \geq 0$. For all $\Lambda^{\text{MMEV}} \in \Delta(\Theta)_{\varepsilon, \delta}^{\text{MMEV}}(R)$, there exists $\pi^{\text{MMEV}} \in$
1159 $\Pi_{\varepsilon, \delta}^{\text{MMEV}}(R)$, such that $(\pi^{\text{MMEV}}, \Lambda^{\text{MMEV}})$ is a (ε, δ) -equilibrium, and we have*

$$V^R(\pi^{\text{MMEV}}; \Lambda^{\text{Deploy}}) = c(\Lambda^{\text{Deploy}}, \alpha - \delta)$$

1160 where $\alpha = V^R(\pi^*; \Lambda^{\text{MMEV}})$ and $c(\Lambda^{\text{Deploy}}, \alpha) = \mathbb{E}_{\theta \sim \Lambda^{\text{Deploy}}} [c(\theta, \alpha)]$

1161 *Proof.* Consider π^* as the optimal policy across all levels $\theta \in \Theta$. At equilibrium, the adversary will
1162 play a strategy

$$\Lambda^{\text{MMEV}} \in \arg\text{-}\delta\text{-}\min_{\Lambda' \in \Delta(\Theta)} V^R(\pi^*; \Lambda')$$

1163 Note it must be the case that $|V^R(\pi^*; \theta) - V^R(\pi^*; \theta')| \leq \delta$ for all θ, θ' in the support of Λ^{MMEV} .
 1164 Also, $\forall \theta \in \Theta$, and θ' in the support of Λ^{MMEV} , we have

$$V^R(\pi^*; \theta) - V^R(\pi^*; \theta') \geq -\delta$$

1165 Considering an arbitrary distribution Λ' , it also holds that

$$V^R(\pi^*; \Lambda') - V^R(\pi^*; \Lambda^{\text{MMEV}}) \geq -\delta$$

1166 In words, the adversary will play levels in which the lowest possible score is achieved by the optimal
 1167 policy.

1168 We now need to construct our policy π^{MMEV} . Take the policy such that

$$V^R(\pi^{\text{MMEV}}; \theta) = \begin{cases} V^R(\pi^*; \theta), & \text{if } \theta \in \text{support } \Lambda^{\text{MMEV}}, \\ c(\theta, \alpha - \delta), & \text{if } \theta \notin \text{support } \Lambda^{\text{MMEV}}. \end{cases}$$

1169 We note that we always have $c(\theta, \alpha - \delta) \geq V^R(\pi^*; \theta')$, where $\theta \in \Theta, \theta' \in \text{support } \Lambda^{\text{MMEV}}$ by
 1170 our definitions. In words, take the MMEV policy such that it achieves the maximum return possible
 1171 on levels played by the adversary, and the smallest available return which is greater than or equal
 1172 to the value achieved on the levels at equilibrium. A policy that achieves this return is clearly in an
 1173 (ε, δ) equilibrium with the adversary we defined at the beginning. The above clearly holds also if
 1174 we consider an arbitrary distribution of levels, i.e. $\exists \pi^{\text{MMEV}}$ such that, for all $\Lambda' \in \Delta(\Theta)$

$$V^R(\pi^{\text{MMEV}}; \Lambda') = c(\Lambda', \alpha - \delta).$$

1175 Taking $\Lambda' = \Lambda^{\text{Deploy}}$ we get

$$V^R(\pi^{\text{MMEV}}; \Lambda^{\text{Deploy}}) = c(\Lambda^{\text{Deploy}}, \alpha - \delta).$$

1176 as desired. □

1177 In our theorem, we proved that an MMEV agent will possibly achieve returns that are at most (close
 1178 to) the ones achievable in levels played by the adversary. If very low return levels exists, this policy
 1179 would then possibly goal misgeneralize at test time if evaluated on levels where a high return is
 1180 possible.

1181 Let's consider an additional small assumptions, i.e. that given a level with highest achievable return
 1182 α_θ , we assume that for all $\theta \in \Theta$ any return in $[0, \alpha_\theta]$ could also be attained.

1183 **Corollary 1.** *Given a level with highest achievable return α_θ , we assume that for all $\theta \in \Theta$ any*
 1184 *return in $[0, \alpha_\theta]$ could also be attained. Additionally consider all the conditions of Theorem 4. Then,*

$$V^R(\pi^{\text{MMEV}}; \Lambda^{\text{Deploy}}) \leq V^R(\pi^*; \Lambda^{\text{MMEV}}) - \delta.$$

1185 *Proof.* See that if every return is attainable (up to the maximum) in each level, we then have
 1186 $c(\Lambda^{\text{Deploy}}, \alpha - \delta) = V^R(\pi^*; \Lambda^{\text{MMEV}}) - \delta$ □

1187 I.2 Training methods and experimental results

1188 In this section, we outline our empirical evaluation of MMEV-based training methods for mitigating
 1189 goal misgeneralization. We adapt three UED methods: PLR^\perp along with two ACCEL variants
 1190 ($\text{ACCEL}^{\text{identity}}$ and $\text{ACCEL}^{\text{binomial}}$, see Appendix H). These UED methods were originally designed
 1191 for MMER training, but we can convert their regret-maximizing adversaries into return-minimizing
 1192 adversaries simply by replacing the regret estimator used to refine the buffer with a negative expected
 1193 return estimator,

$$\hat{M}_{\text{latest}}^R(\pi; \theta) = -\hat{V}_{\text{latest}}^R(\pi; \theta), \quad (16)$$

1194 where $\hat{V}_{\text{latest}}^R(\pi; \theta)$ is the empirical average return achieved on θ in the latest batch of rollouts with
 1195 the current policy π .

1196 Figure I.1 show the performance of trained policies in CHEESE IN THE CORNER and KEYS AND
 1197 CHESTS for ambiguous and distinguishing levels. As expected, these training methods fail to miti-
 1198 gate goal misgeneralization, and even fail to induce robust performance in ambiguous levels.

1199 We observe that these adversaries rapidly fill their buffers with levels with zero estimated value,
 1200 indicating that the adversaries are working as expected. We hypothesize that the training failure
 1201 is primarily due to the adversary finding levels in which the policy not only receives low expected
 1202 value, but never receives any nonzero reward and thus obtains no training signal.

1203 In the extreme case, the adversary could populate the buffer with levels that have zero *maximum*
 1204 value, preventing any policy from obtaining nonzero reward. In CHEESE IN THE CORNER, levels
 1205 in which the cheese position is unreachable from the mouse spawn position have zero maximum
 1206 value. In Figure I.2, we plot the average proportion of such “unsolvable” levels in the adversary’s
 1207 buffer over training for CHEESE IN THE CORNER. We find that this proportion is around 80% for
 1208 most training distributions, which is much higher than the baseline value of around 18% of levels
 1209 sampled from $\Lambda_{\text{Ambig.}}$ that are unsolvable. As the proportion of distinguishing levels increases, the
 1210 average proportion of unsolvable levels decreases slightly, to around 50% for $\alpha = 0$. Note that the
 1211 baseline value for $\Lambda_{\text{Distg.}}$ is also lower (around 7%) since it’s easier for walls to obstruct the cheese
 1212 when it’s in the corner than when it is in an arbitrary position in the interior of the grid. Therefore we
 1213 hypothesize that the adversaries have a slightly harder time finding unsolvable levels as α increases.
 1214 However, unsolvable levels still occupy a majority of the buffer.

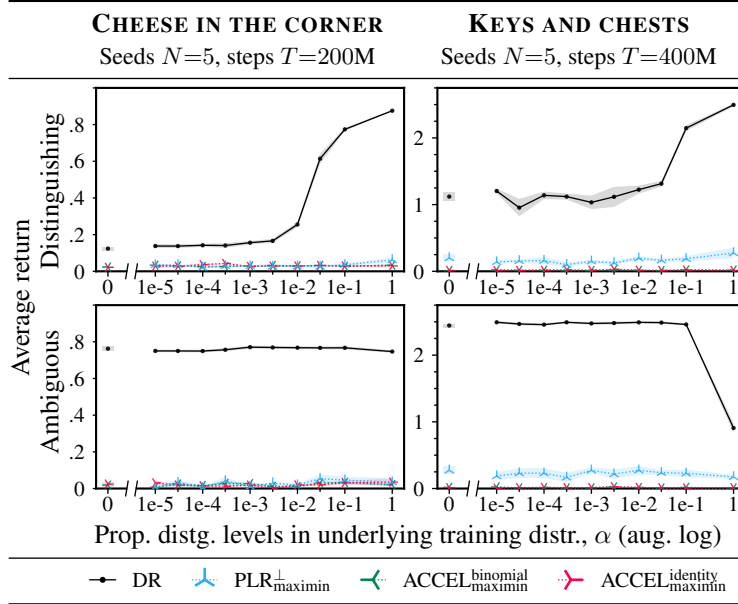


Figure I.1: Maximin training methods. Each policy is trained on T environment steps, using the indicated algorithm, with underlying training distribution $\Lambda_{\alpha}^{\text{Train}}$. Average return over 512 steps for an evaluation batch of 256 levels (top row: distinguishing levels, bottom row: ambiguous levels). Mean over N seeds, shaded region is one standard error. Note the split axes used to show zero on the log scale. Note the drop in performance for DR in ambiguous KEYS AND CHESTS levels at $\alpha = 1$ can be explained by noting that ambiguous levels are now out of distribution given this training distribution.

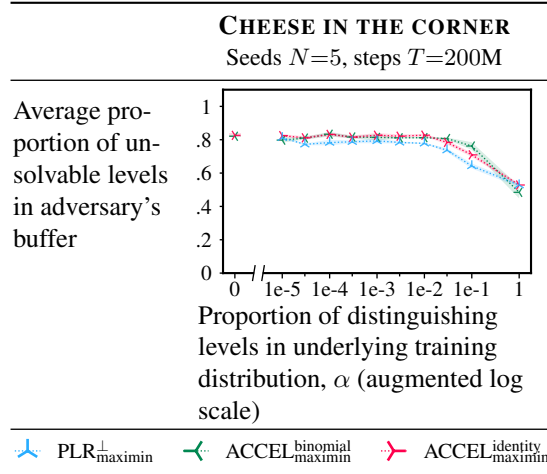


Figure I.2: Maximin training methods. Each policy is trained on T environment steps, using the indicated algorithm, with underlying training distribution $\Lambda_{\alpha}^{\text{Train}}$. Average proportion of unsolvable levels in the adversary's buffer over training. Mean over N seeds, shaded region is one standard error. Note the split axes used to show zero on the log scale.

1215 J Experiments with increased training length

1216 In most of our experiments, we compare the performance of algorithms after a fixed amount of
 1217 training time, and find that MMER-based training methods typically outperform MEV-based train-
 1218 ing methods for a fixed training budget. We have shown that this is because a regret-maximizing
 1219 adversary can amplify the proportion of distinguishing levels compared to sampling from a fixed
 1220 underlying training distribution.

1221 Another method of increasing the agent’s experience in distinguishing levels is to train for longer.
 1222 In this appendix, we extend training times in the CHEESE ON A DISH training environment (where
 1223 DR was most robust to goal misgeneralization).

1224 Figure J.1 shows the results. We find that training for more than 200 million environment steps with
 1225 a fixed training method slightly mitigates goal misgeneralization. In particular, for $\alpha = 1e-3$, DR
 1226 gradually stops misgeneralizing after 200 million steps. The result is qualitatively with Theorem 1,
 1227 since increasing training time should have the effect of decreasing the optimization threshold.

1228 However, for lower α values, further training shows diminishing returns, and even 1,500 million
 1229 steps of DR training is insufficient to mitigate goal misgeneralization to the extent achieved by most
 1230 UED methods within 200 million steps. This suggests that current UED methods are both more effi-
 1231 cient and also qualitatively more effective at mitigating goal misgeneralization in this environment.

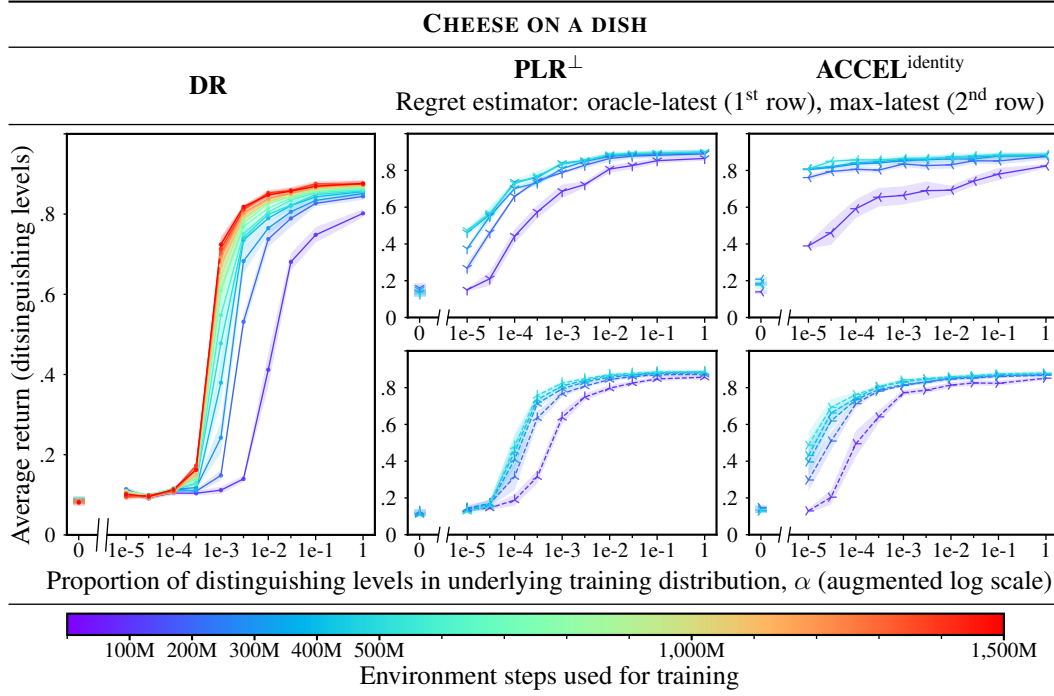


Figure J.1: Training for more environment steps in CHEESE ON A DISH. Each policy is trained for 500M environment steps (1,500M for DR), using the indicated algorithm, with underlying training distribution $\Lambda_{\alpha}^{\text{Train}}$. Every 100M steps, we evaluate the average return over 512 steps for an evaluation batch of 256 distinguishing levels (cf. Figure 3). Mean over 3 seeds, shaded region is one standard error. Note the split axes used to show zero on the log scale.

1232 K Experiments with a different distinguishing level generator

1233 In this section, we consider an alternative procedural level generator for distinguishing levels in the
 1234 CHEESE IN THE CORNER environment. Our main experiments use a distinguishing level generator
 1235 that places the cheese anywhere in the maze.

1236 We consider a restricted distinguishing level generator that positions the cheese only within a region
 1237 of size $c \times c$ surrounding the top-left corner, for varying c (the ambiguous generator would be re-
 1238 covered with $c = 1$, and the original, unrestricted distinguishing level generator would be recovered
 1239 with $c = 13$).

1240 Figure K.1 shows the return is evaluated on unrestricted distinguishing levels, where the cheese is
 1241 positioned anywhere in the maze. We find that the UED methods are able to amplify the proportion
 1242 of restricted distinguishing levels and in some cases mitigate goal misgeneralization. DR achieves
 1243 low return across all corner sizes c .

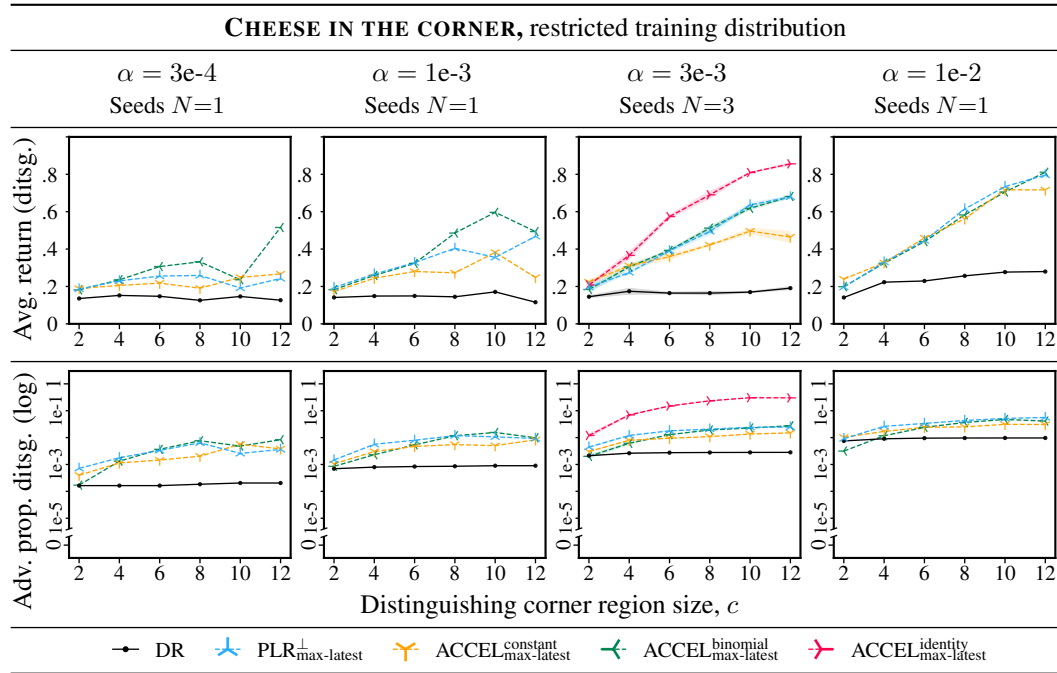


Figure K.1: **Training with varying distinguishing level generators.** Each policy is trained on 200 million environment steps, using the indicated algorithm, with underlying training distribution $\Lambda_{\alpha,c}^{\text{Train}} = (1 - \alpha)\Lambda_{\text{Ambig.}} + \alpha\Lambda_{\text{Distg.}}^c$, where $\Lambda_{\text{Distg.}}^c$ is a procedural level generator that positions the cheese in the top-left $c \times c$ region of the maze. (1st row): Average return over 512 steps for an evaluation batch of 256 distinguishing levels (cf. Figure 3). (2nd row): The proportion of distinguishing levels sampled from the adversary’s buffer across training (cf. Figure 4). (Both): Mean over N seeds, shaded region is one standard error. Note the split vertical axis used to show zero on the log scale.

1244 L Experiments with different observations

1245 In this section, we vary the way we encode observations for the policy in CHEESE ON A DISH and
 1246 explore its effect on goal misgeneralization.

1247 As discussed in Appendix D.2, the proxy goal and the true goal are symmetric in this environment,
 1248 other than the fact that the position of the dish is redundantly represented across multiple channels in
 1249 the Boolean grid observation. In the main text, we use $D = 6$ channels to encode the dish position
 1250 (compared to 1 channel for the cheese position). The additional channels break a symmetry and
 1251 create a slight inductive bias in favor of a policy that pursues the proxy goal.

1252 We conduct an experiment where we vary the number of channels and see what affect it has on goal
 1253 misgeneralization. Figure L.1 shows the results. We see that with $D = 1$ channel coding the dish
 1254 position, all methods (including DR) are somewhat robust to goal misgeneralization, even at small
 1255 α values. Additional channels significantly increase DR’s susceptibility to goal misgeneralization.
 1256 On the other hand, UED methods remain able to identify and amplify the training signal from rare
 1257 distinguishing levels, while UED methods retain comparably similar performance.

1258 The extent of amplification remains essentially constant with D . We hypothesize that the number of
 1259 channels does not stop the adversary from noticing high-regret distinguishing levels, though it may
 1260 affect how the policy responds.

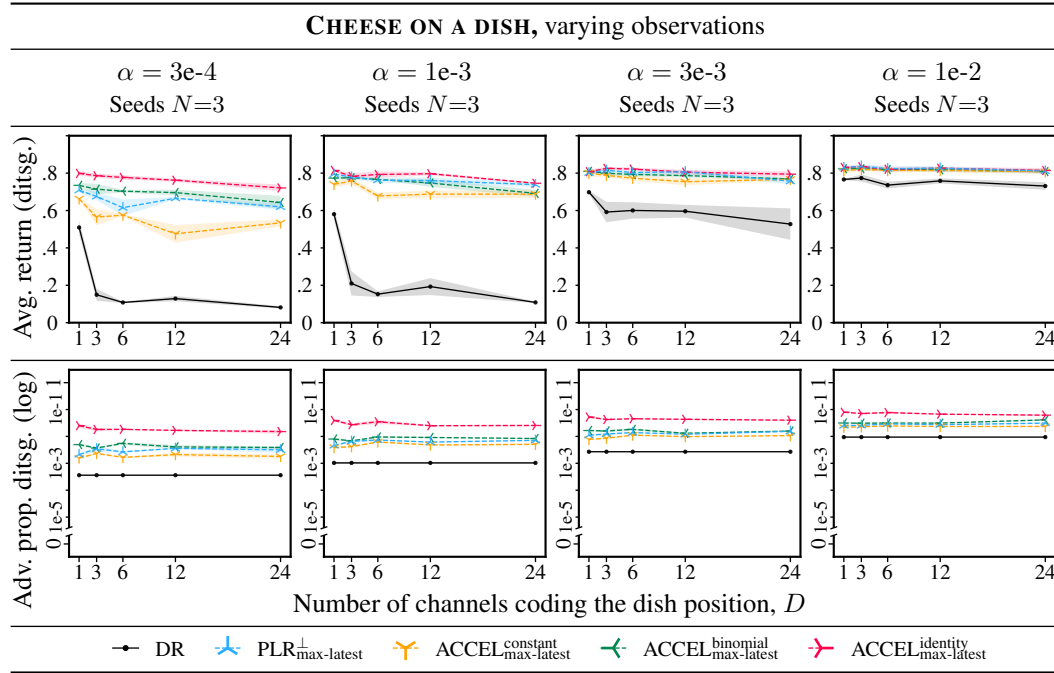


Figure L.1: **Training with observations with varying emphasis on the dish.** Each policy is trained on 200 million environment steps, using the indicated algorithm, with underlying training distribution $\Lambda_{\alpha}^{\text{Train}}$. We vary the number of channels (features) encoding the dish position, D . (1st row): Average return over 512 steps for an evaluation batch of 256 distinguishing levels (cf. Figure 3). (2nd row): The proportion of distinguishing levels sampled from the adversary’s buffer across training (cf. Figure 4). (Both): Mean over N seeds, shaded region is one standard error. Note the split vertical axis used to show zero on the log scale.