# Mitigating Goal Misgeneralization via Minimax Regret

**Brought together by KASL** — KRUEGER AI SAFETY LAB

**Karim Abdel Sadek(=)**
University of California, Berkeley

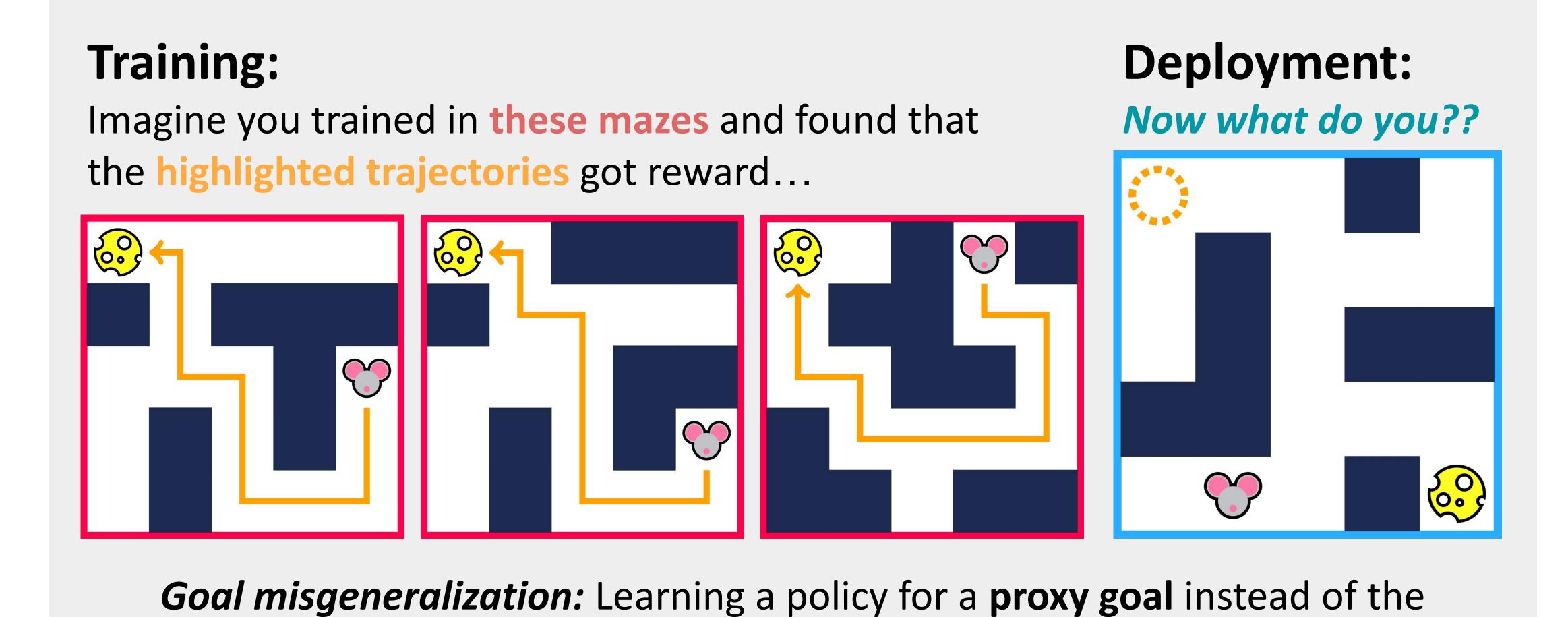**Matthew Farrugia-Roberts(=)**
University of Oxford

**Usman Anwar**
University of Cambridge

**Hannah Erlebach**
University College London

**Christian Schroeder de Witt**
University of Oxford
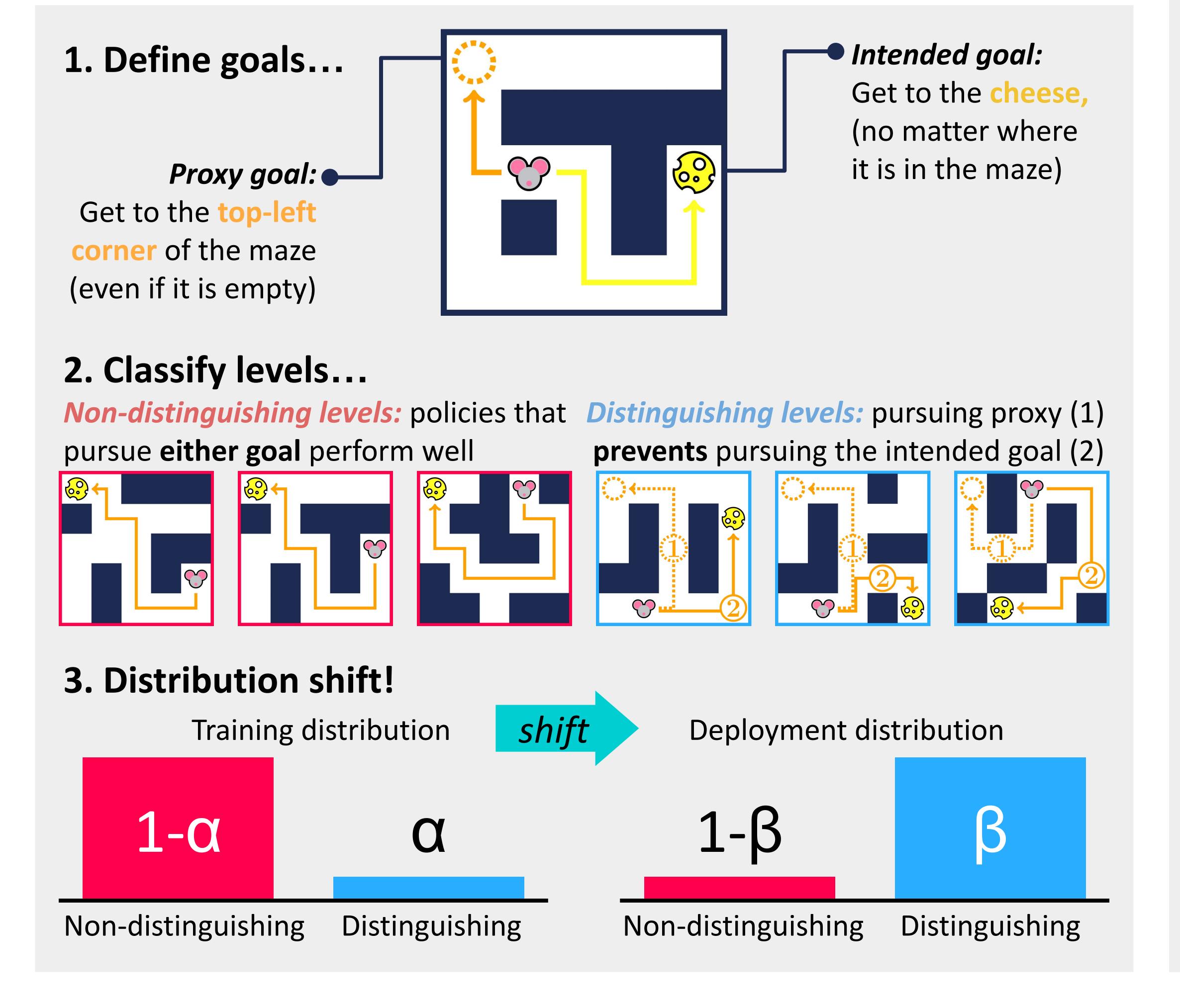
**David Krueger**
Mila, University of Montreal

**Michael Dennis**
Google DeepMind

## Quiz: Are *YOU* susceptible to *goal misgeneralization??*

**Training:**
Imagine you trained in **these mazes** and found that the **highlighted trajectories** got reward…

**Deployment:**
*Now what do you??*



*Goal misgeneralization:* Learning a policy for a **proxy goal** instead of the **intended goal** from an **ambiguous training environment** distribution.

## Problem Setting:
*Proxy-Distinguishing Distribution Shift*

**1. Define goals…**

*Proxy goal:*
Get to the **top-left corner** of the maze (even if it is empty)

*Intended goal:*
Get to the **cheese,** (no matter where it is in the maze)

**2. Classify levels…**

*Non-distinguishing levels:* policies that pursue **either goal** perform well

*Distinguishing levels:* pursuing proxy (1) **prevents** pursuing the intended goal (2)

**3. Distribution shift!**

Training distribution — **shift** → Deployment distribution

| Non-distinguishing | Distinguishing |
| --- | --- |
| 1-α | α |

| Non-distinguishing | Distinguishing |
| --- | --- |
| 1-β | β |

---

We show that training with **the *maximum expected value* objective** is **susceptible** to goal misgeneralization!

Approximate **Maximum Expected Value (MEV)** objective:

$$\pi^{\text{MEV}} \in \underset{\pi \in \Pi}{\text{arg-}\varepsilon\text{-max}} \, \text{Value}(\pi, \Lambda^{\text{Train}})$$

**Approx. maximization**
within some threshold $\varepsilon \geq 0$ of optimal policy

**Expected return**
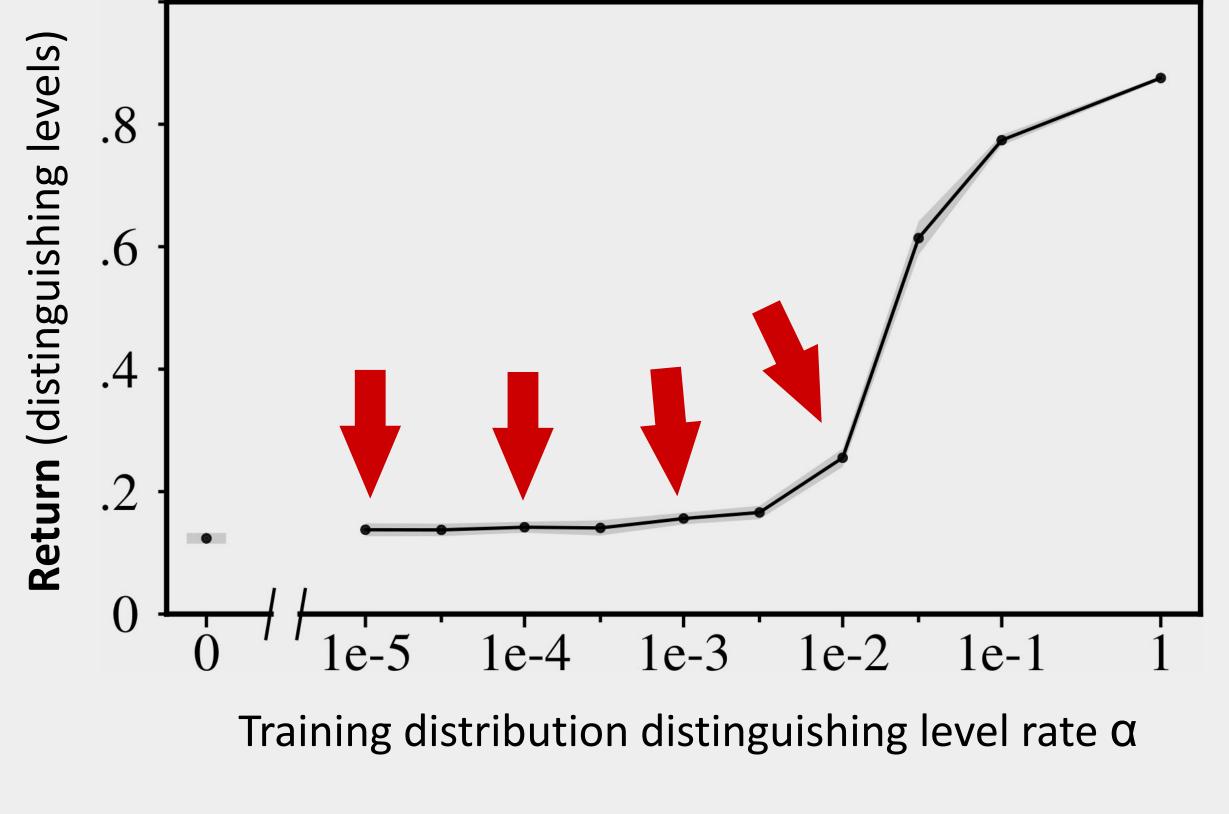over some fixed level distribution

**Training distribution**
of non-distinguishing / distinguishing levels

**Theorem 1:** If $\alpha \leq \varepsilon$, some MEV policies pursue the proxy goal:

$$\exists \pi^{\text{MEV}}; \pi^{\text{MEV}} \in \underset{\pi \in \Pi}{\arg\max} \, \text{ProxyValue}(\pi, \Lambda^{\text{Deploy}}) \setminus \underset{\pi \in \Pi}{\text{arg-}\beta\text{-max}} \, \text{Value}(\pi, \Lambda^{\text{Deploy}})$$

**Experiments with *Domain Randomization:***
We train with **domain randomization** (implementing the MEV objective). We use training distributions with **varying α** (proportion of distinguishing levels).



When α < 0.03, domain randzn. learns a policy that **fails to pursue the intended goal** on distinguishing levels...

… instead the policy **pursues the proxy goal** on these levels, leading to **misgeneralization in deployment.**

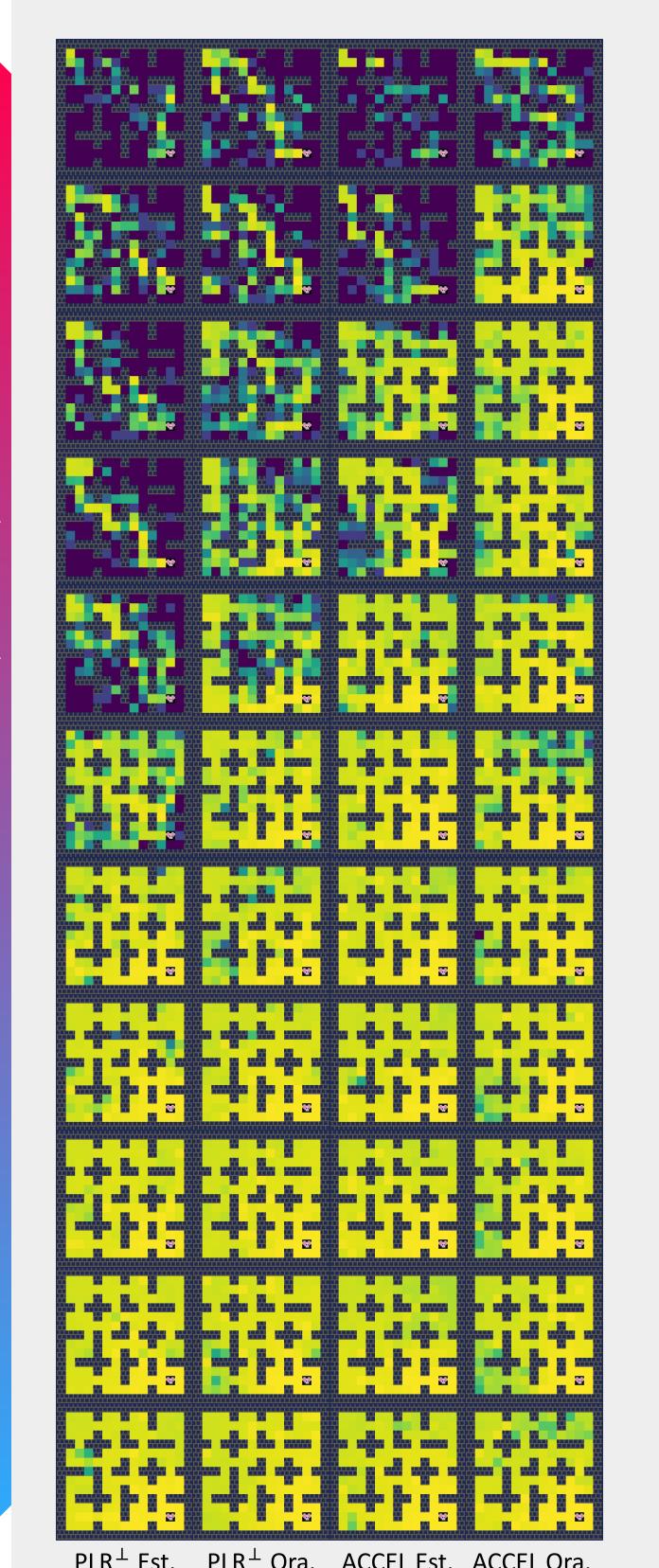*Blind spots!*
For each *position*

*Return* when cheese is there

---

On the other hand, training with **the *minimax expected regret* objective** is **robust** to goal misgeneralization!

Approximate **MiniMax Expected Regret (MMER)** objective:

$$\pi^{\text{MMER}} \in \underset{\pi \in \Pi}{\text{arg-}\varepsilon\text{-min}} \, \underset{\Lambda \in \Delta(\text{lvl.})}{\max} \, \text{Regret}(\pi, \Lambda)$$

**Approx. *minimization***
within some threshold $\varepsilon \geq 0$ of optimal policy

***Inner maximization***
worst-case level distr. relative to policy

**Expected *Regret*:**
$\text{Value}(\pi^{\star}, \text{level}) - \text{Value}(\pi, \text{level})$
averaged over level distribution

**Theorem 2:** All MMER policies pursue the intended goal:

$$\forall \pi^{\text{MMER}}, \pi^{\text{MMER}} \in \underset{\pi \in \Pi}{\text{arg-}\varepsilon\text{-max}} \, \text{Value}(\pi, \Lambda^{\text{Deploy}})$$

**Experiments with *Unsupervised Environment Design:***
Train with **unsupervised environment design** (implementing MMER objective). We use four increasingly powerful adversarial designers and regret estimators.



UED policies **pursue the intended goal** at many low α where domain randomization policy pursued the proxy goal.

**How?** The adversary **finds high-regret distinguishing levels** and plays them more often than α.

*See paper for…*
+ theory details
+ more results
+ more environments
+ more methods