

# Universal Adaptive Environment Discovery

Madi Matymov\*  
KAUST  
Saudi Arabia

Ba-Hien Tran  
Huawei Technologies France SASU  
France

Maurizio Filippone  
KAUST  
Saudi Arabia

## Abstract

An open problem in Machine Learning is how to avoid models to exploit spurious correlations in the data; a famous example is the background-label shortcut in the Waterbirds dataset. A common remedy is to train a model across multiple *environments*; in the Waterbirds dataset, this corresponds to training by randomizing the background. However, selecting the right environments is a challenging problem, given that these are rarely known a priori. We propose *Universal Adaptive Environment Discovery* (UAED), a unified framework that *learns* a distribution over data transformations that instantiate environments, and optimizes any robust objective *averaged* over this learned distribution. UAED yields adaptive variants of IRM, REx, GroupDRO, and CORAL without predefined groups or manual environment design. We provide a theoretical analysis by providing PAC-Bayes bounds and by showing robustness to test environment distributions under standard conditions. Empirically, UAED discovers interpretable environment distributions and improves worst-case accuracy on standard benchmarks, while remaining competitive on mean accuracy. Our results indicate that making environments *adaptive* is a practical route to out-of-distribution generalization.

## 1 Introduction

Machine learning models often fail when deployed in environments that differ from their training conditions. For example, a medical diagnosis system trained on data from one hospital may perform poorly at another, and a model trained to recognize birds from

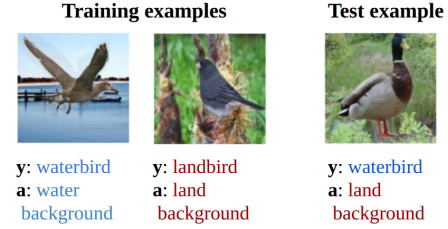


Figure 1: waterBirds dataset examples (Sagawa et al., 2020b): the training correlation between label  $y$  and spurious attribute  $a$  does not hold at test time.

professional photographs may struggle with amateur images. Such failures typically arise because models rely on *spurious correlations* (see, e.g., Sagawa et al. (2020a), and references therein)—patterns present in the training data that do not generalize to deployment settings. Fig. 1 illustrates this phenomenon on the WaterBirds dataset (Sagawa et al., 2020b). When a model learns spurious correlations, it can achieve high i.i.d. test accuracy but fail dramatically on subgroups where the correlation does not hold.

The machine learning community has developed numerous approaches to address this challenge. Invariant Risk Minimization (Arjovsky et al., 2019) seeks features with stable predictive relationships across environments. Risk Extrapolation (REX) (Krueger et al., 2021) minimizes variance of risks across environments. Group Distributionally Robust Optimization (GroupDRO) (Sagawa et al., 2020a) optimizes worst-case performance over predefined groups. Domain alignment methods include covariance alignment via CORAL, which matches second-order feature statistics across domains (Sun and Saenko, 2016), and kernel mean-matching approaches based on Maximum Mean Discrepancy (MMD) (Gretton et al., 2012).

Despite their different objectives, all these methods share a critical limitation: they require practitioners to predefine environments or groups that appropriately expose spurious correlations. This creates a fundamental paradox—identifying which correlations are spurious requires knowing the environments, but defining good environments requires knowing which corre-

\*Corresponding author: madi.matymov@kaust.edu.sa

lations are spurious.

**Our Contributions.** We introduce Universal Adaptive Environment Discovery (UAED), a framework that makes environment specification learnable rather than fixed. Specifically, we: **(1)** *unify diverse robust methods* (IRM, REX, GroupDRO, CORAL) under adaptive environment discovery, parameterized through learnable data transformations; **(2)** *provide a theoretical foundation* via PAC-Bayes bounds showing distributionally robust guarantees with implicit regularization explaining method-specific behaviors; **(3)** *validate empirically* on synthetic and real-world benchmarks, where adaptive variants consistently improve over baselines and approach state-of-the-art without predefined groups; and **(4)** *offer the conceptual insight* that robust objectives and environment discovery are complementary, with the latter being critical to the former’s success.

Our work reveals that the dichotomy between “choosing robust objectives” and “defining environments” is artificial. By making environments learnable, we provide a more principled and practical path toward out-of-distribution generalization.

## 2 Related Work

**Robust Learning Methods.** There exist several approaches for robust learning, each with distinct objectives. IRM (Arjovsky et al., 2019) seeks invariant predictors but is highly sensitive to environment specification (Gulrajani and Lopez-Paz, 2021; Rosenfeld et al., 2021). Variants such as IB-IRM (Ahuja et al., 2020) and PAIR (Kim et al., 2021) still require predefined environments. REX (Krueger et al., 2021) and Risk Variance Penalization (RVP) (Xie et al., 2020) minimize variance across environments, while GroupDRO (Sagawa et al., 2020a) optimizes worst-case group performance. More recent methods, including Just Train Twice (JTT) (Liu et al., 2021a), DISC (Li et al., 2023), and Deep Feature Reweighting (DFR) (Kirichenko et al., 2023), achieve strong performance but rely on group annotations or specialized architectures.

**Environment Discovery.** Prior work has approached environment discovery from different perspectives. EIIL (Creager et al., 2021) leverages causal discovery but requires access to the underlying causal graph. Environment inference methods such as heterogeneous risk minimization (Liu et al., 2021b) generally assume environments are discrete. LISA (Yao et al., 2022) instead employs MIXUP (Zhang et al., 2018) for implicit environment augmentation. In contrast, our framework unifies these directions by enabling any robust objective to operate with learned *continuous* environment distributions.

**Data Augmentation Learning.** AUTOAUGMENT (Cubuk et al., 2019) and RANDAUGMENT (Cubuk et al., 2020) learn augmentation policies for standard training, while AUGERINO (Benton et al., 2020) learns invariances directly from data. AUGMAX (Wang et al., 2019) employs adversarial augmentations to enhance robustness. Most relevant is targeted augmentation (Gao et al., 2023), which shows that augmentation can improve OOD performance. OPTIMA (Matymov et al., 2025) placed priors over transformation policies and performs evidence-driven selection to improve generalization under shift. We build on these insights by showing that augmentation learning and environment specification are fundamentally connected.

**Theoretical Foundations.** PAC-Bayes theory (McAllester, 1999) provides generalization bounds that trade off empirical risk against model complexity. For domain adaptation, Germain et al. (2013) derived bounds based on divergence measures. More recently, Deng et al. (2020) established connections between PAC-Bayes and domain generalization. We extend this line of work by proving that adaptive environment discovery yields distributionally robust guarantees for any robust learning objective.

## 3 Universal Adaptive Environment Discovery

**Setup.** We consider inputs  $X \in \mathcal{X}$ , labels/targets  $Y \in \mathcal{Y}$ , and environments  $e \in \mathcal{E}$ . Each environment induces a joint distribution  $P^e$  on  $\mathcal{X} \times \mathcal{Y}$ . The goal is to learn a predictor  $f_\theta : \mathcal{X} \rightarrow \hat{\mathcal{Y}}$  that *generalizes across environments*, even when the marginals  $P^e(X)$  and spurious dependencies in  $P^e(Y | X)$  vary with  $e$ . We assume the existence of an underlying stable mechanism, formalized next.

**Assumption 3.1** (Invariant conditional). There exists a representation  $\Phi^* : \mathcal{X} \rightarrow \mathcal{Z}$  such that the conditional distribution of  $Y$  given the representation is invariant across environments:

$$P^e(Y | \Phi^*(X)) = P(Y | \Phi^*(X)) \quad \text{for all } e \in \mathcal{E}.$$

**Remarks.** (i) This is the standard invariant-prediction premise (e.g., (Peters et al., 2016; Arjovsky et al., 2019)); it is *task-agnostic* and covers classification and regression (no specific link function is assumed). (ii) We use it to motivate environment diversity; our PAC-Bayes and DRO guarantees do *not* rely on Assumption 3.1.

### 3.1 Baseline robust objectives (fixed environments)

Let  $\mathcal{E}_{\text{train}} = \{e_1, \dots, e_k\}$  be given. For a bounded or sub-gamma loss  $\ell$ ,

$$\mathcal{R}^e(\theta) = \mathbb{E}_{(x,y) \sim P^e} \ell(f_\theta(x), y).$$

We write all baselines (except GroupDRO) as

$$\min_{\theta} \underbrace{\frac{1}{k} \sum_{e \in \mathcal{E}_{\text{train}}} \mathcal{R}^e(\theta)}_{\text{mean risk}} + \eta \underbrace{\mathcal{P}_{\text{robust}}^{\text{fixed}}(\theta; \mathcal{E}_{\text{train}})}_{\text{method-specific regularizer}}. \quad (1)$$

**IRM (v1)** (Arjovsky et al., 2019):

$$\mathcal{P}_{\text{robust}}^{\text{fixed}} = \frac{1}{k} \sum_e \left\| \nabla_w \mathcal{R}^e(w \cdot f_\theta) \Big|_{w=1} \right\|_2^2.$$

**REx / VREx** (Krueger et al., 2021):

$$\mathcal{P}_{\text{robust}}^{\text{fixed}} = \text{Var}_e [\mathcal{R}^e(\theta)].$$

**CORAL** (Sun and Saenko, 2016): Let  $F_\theta(\cdot)$  denote features. Using minibatch covariances,

$$\mathcal{P}_{\text{robust}}^{\text{fixed}} = \frac{1}{k(k-1)} \sum_{e \neq e'} \left\| \text{Cov}(F_\theta(X^e)) - \text{Cov}(F_\theta(X^{e'})) \right\|_F^2.$$

**GroupDRO** (Sagawa et al., 2020a) (separate form):

$$\min_{\theta} \max_{e \in \mathcal{E}_{\text{train}}} \mathcal{R}^e(\theta),$$

often optimized via the entropic surrogate  $\min_{\theta} \frac{1}{\lambda} \log \left( \frac{1}{k} \sum_e e^{\lambda \mathcal{R}^e(\theta)} \right)$  with  $\lambda > 0$ .

All these approaches assume access to predefined environments  $\{e_1, \dots, e_k\}$ , leading to the central challenge: *how can we define environments without knowing a priori what distinguishes them?*

**Policy over environments.** Instead of fixing  $\{e_1, \dots, e_k\}$ , we index environments by a parameter  $\gamma \in \Gamma$  (e.g., correlation strength, rotation angle, style). A *policy*  $\Pi_\phi = p(\gamma \mid \phi)$  (density or categorical mass) over  $\Gamma$  selects environments during training; let  $e(\gamma)$  be the corresponding environment.

### 3.2 The Unified UAED Framework

Inspired by variational principles (Jordan et al., 1999), we replace fixed environments with a learned distribution over environments. We define the policy-averaged risk as

$$\mathcal{R}_{\Pi_\phi}(\theta) = \mathbb{E}_{\gamma \sim \Pi_\phi} \mathbb{E}_{(x,y) \sim P} \ell(f_\theta(T_\gamma(x)), y),$$

where  $T_\gamma$  denotes the (possibly stochastic) data transformation that instantiates  $e(\gamma)$ .

**Definition 3.2** (Universal Adaptive Objective). Given a robust method with regularizer  $\mathcal{P}_{\text{robust}}$ , UAED optimizes

$$\min_{\theta, \phi} \underbrace{\mathbb{E}_{\gamma \sim \Pi_\phi} [\mathcal{R}^{e(\gamma)}(\theta)]}_{\text{policy-averaged risk}} + \underbrace{\eta \mathcal{P}_{\text{robust}}(\theta; \Pi_\phi)}_{\text{method-specific}} + \beta \text{KL}(\Pi_\phi \parallel \Pi_0), \quad (2)$$

where  $\Pi_0$  is a fixed prior on  $\Gamma$ , and  $\eta, \beta \geq 0$ .

For each method,  $\mathcal{P}_{\text{robust}}(\theta; \Pi_\phi)$  is

**Adaptive IRM (A-IRM):**

$$\mathbb{E}_{\gamma \sim \Pi_\phi} \left\| \nabla_w \mathcal{R}^{e(\gamma)}(w \cdot f_\theta) \Big|_{w=1} \right\|_2^2.$$

**Adaptive REx (A-REx):**

$$\text{Var}_{\gamma \sim \Pi_\phi} [\mathcal{R}^{e(\gamma)}(\theta)].$$

**Adaptive CORAL (A-CORAL):**

$$\mathbb{E}_{\gamma_1, \gamma_2 \sim \Pi_\phi} \left\| \text{Cov}(F_\theta(X^{e(\gamma_1)})) - \text{Cov}(F_\theta(X^{e(\gamma_2)})) \right\|_F^2.$$

**Adaptive GroupDRO (A-GroupDRO):** Directly maximizing over  $\gamma$  is non-smooth. Using the entropic-risk dual of a KL-ball DRO (Prop. 4.5), we obtain a smooth surrogate:

$$\min_{\theta, \phi} \underbrace{\frac{1}{\lambda} \log \mathbb{E}_{\gamma \sim \Pi_\phi} \exp \{ \lambda \mathcal{R}^{e(\gamma)}(\theta) \}}_{\text{entropic (smooth) worst-case}} + \beta \text{KL}(\Pi_\phi \parallel \Pi_0),$$

for  $\lambda > 0$ , optionally adding  $\rho/\lambda$  to enforce robustness to a KL-ball of radius  $\rho$  around  $\Pi_\phi$ .

**Variational note.** When  $\ell$  is a negative log-likelihood, the UAED objective in (2) is (up to the method-specific regularizer and a temperature  $\beta$  (Higgins et al., 2017; Bissiri et al., 2016)) exactly a negative ELBO with prior  $\Pi_0$  over  $\gamma$  and variational posterior  $\Pi_\phi$ ; our Loss Averaging (LoA)-Log-Likelihood Averaging (LA) Theorem 4.6 bound justifies using the policy-averaged risk in place of the log-mixture.

### 3.3 Implementation: Hierarchical Bayesian Framework

We adopt a hierarchical Bayesian framework (see, e.g., Gelman et al., 1995) for the environment policy  $p(\gamma \mid \phi)$  to encourage broad exploration and mitigate the risk of overfitting to specific, narrow environment.

**Continuous Environments:** For continuous  $\gamma \in [0, 1]$  (e.g., correlation strength), we model the policy

as a Beta distribution whose shape and rate parameters are generated by a multi-layer perceptron (MLP):

$$\gamma \sim \text{Beta}(\alpha(\phi), \beta(\phi)), \quad (3)$$

$$\alpha(\phi), \beta(\phi) = \text{softplus}(\text{MLP}(\phi)) + \varepsilon, \quad (4)$$

where  $\varepsilon = 10^{-6}$  is a small *deterministic* offset used only for numerical stability.

**Discrete Environments:** For discrete transformations, the environment is modeled by a Categorical distribution with logits parameterized by  $\phi$ :

$$p(\gamma|\phi) = \text{Categorical}(\text{softmax}(\phi/\tau)). \quad (5)$$

To enable end-to-end backpropagation, we employ Gumbel-Softmax reparameterization (Jang et al., 2017) to sample from this discrete distribution. We employ a temperature schedule  $\tau_t$  annealed linearly from 1.0 to 0.3 over the first  $T_0$  epochs.

### 3.4 Why Different Methods Need Different Environments

Each robust objective implicitly seeks different environmental properties: A-IRM discovers environments with conflicting spurious correlations, A-REx favors maximally diverse yet learnable environments, A-GroupDRO identifies stress-test scenarios, and A-CORAL aligns environments with distinct marginals. This flexibility to tailor the discovery process is the core advantage of our approach, enabling each method to automatically generate the environments needed for its objective to converge to the true invariant mechanism.

## 4 Theoretical Analysis

In this section, we give guarantees for UAED that (i) are PAC-Bayes for the loss averaged under the *learned* environment policy  $\Pi_\phi$ , and (ii) imply *distributional robustness in environment space*: any test environment distribution inside a KL-ball around  $\Pi_\phi$  is controlled.

**Setup.** Let  $X \in \mathcal{X}$ ,  $Y \in \mathcal{Y}$ ,  $Z = (X, Y) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ , and  $P$  a base distribution on  $\mathcal{Z}$ . Let  $\Gamma$  index data transformations and  $\Pi_\phi$  be a policy over  $\Gamma$ . Each  $\gamma \in \Gamma$  specifies a (possibly stochastic) map  $T_\gamma$  that maps a point  $z \in \mathcal{Z}$  to a distribution over  $z' \in \mathcal{Z}$ . The induced environment is the *pushforward* of  $P$  by  $T_\gamma$ :

$$P^{e(\gamma)} := (T_\gamma)_\# P, \quad \text{meaning} \quad (6)$$

$$\int f(z') dP^{e(\gamma)}(z') = \int \left( \mathbb{E}_{z' \sim T_\gamma(\cdot|z)} f(z') \right) dP(z)$$

for all bounded measurable  $f$ . For a bounded loss  $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow [0, 1]$  and predictor  $h \in \mathcal{H}$ ,

$$\mathcal{R}^{e(\gamma)}(h) = \mathbb{E}_{z \sim P} \mathbb{E}_{z' \sim T_\gamma(\cdot|z)} \ell(h, z'), \quad (7)$$

$$\mathcal{R}_{\Pi_\phi}(h) = \mathbb{E}_{\gamma \sim \Pi_\phi} \mathcal{R}^{e(\gamma)}(h). \quad (8)$$

Given  $S = (z_1, \dots, z_n) \sim P^{\otimes n}$ , the empirical risk is

$$\hat{\mathcal{R}}_{\Pi_\phi}(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\gamma \sim \Pi_\phi} \mathbb{E}_{z' \sim T_\gamma(\cdot|z_i)} \ell(h, z'). \quad (9)$$

*Readability note.* One can read  $(T_\gamma)_\# P$  as “the distribution of the transformed sample  $Z'$ , obtained by applying  $\gamma$  to  $Z \sim P$ .”

**Joint prior and posterior.** Treat the *model* and the *policy* as a single hypothesis  $H = (h, \phi)$  with prior  $M = P \times \Pi_0$  that is independent of  $S$ . Let  $Q$  be any posterior over  $(h, \phi)$  learned from  $S$ . All expectations  $\mathbb{E}_{H \sim Q}[\cdot]$  are w.r.t. this joint posterior.

**Assumption 4.1** (Bounded or sub-gamma loss). Either (a)  $\ell \in [0, 1]$ , or (b)  $\ell$  is sub-gamma under all  $P^{e(\gamma)}$  with variance proxy  $\sigma^2$  and scale  $c > 0$  (Remark 4.7).

### 4.1 PAC-Bayes under a learned policy

**Theorem 4.2** (PAC-Bayes for policy-averaged risk). *Under Assumption 4.1 with bounded loss, with probability at least  $1 - \delta$  over  $S$ , for all posteriors  $Q$ ,*

$$\mathbb{E}_{H \sim Q} [\mathcal{R}_{\Pi_\phi}(h)] \leq \mathbb{E}_{H \sim Q} [\hat{\mathcal{R}}_{\Pi_\phi}(h)] + \sqrt{\frac{KL(Q||M) + \ln(1/\delta)}{2n}}. \quad (10)$$

The theorem establishes a direct generalization guarantee for the  $\Pi_\phi$ -averaged risk by applying the standard PAC-Bayes bound (McAllester, 1999) to the composite loss  $(h, \phi, z) \mapsto \mathbb{E}_{\gamma \sim \Pi_\phi} \mathbb{E}_{z' \sim T_\gamma(\cdot|z)} \ell(h, z') \in [0, 1]$ . Since the policy  $\phi$  is learned simultaneously with the model  $h$ , the policy-averaged empirical risk  $\mathbb{E}_{H \sim Q} [\hat{\mathcal{R}}_{\Pi_\phi}(h)]$  serves as a sharp proxy for the true policy-averaged risk  $\mathbb{E}_{H \sim Q} [\mathcal{R}_{\Pi_\phi}(h)]$ , controlled by the complexity term  $KL(Q||M)$ .

### 4.2 Robustness to environment shift via KL-balls

We now relate risks under any *test* environment distribution  $G$  over  $\gamma$  to those under  $\Pi_\phi$ .

**Lemma 4.3** (Change of environment via DV+Hoeffding). *Let  $f : \Gamma \rightarrow [0, 1]$  be measurable. For any  $G$  and any policy  $\Pi_\phi$ ,*

$$\mathbb{E}_{\gamma \sim G} f(\gamma) \leq \mathbb{E}_{\gamma \sim \Pi_\phi} f(\gamma) + \sqrt{\frac{1}{2} KL(G||\Pi_\phi)}.$$

**Theorem 4.4** (UAED robust generalization). *Under Assumption 4.1 with bounded loss, with probability at least  $1 - \delta$  over  $S$ , for all posteriors  $Q$  and all  $G$  satisfying  $KL(G||\Pi_\phi) \leq \rho$ ,*

$$\mathbb{E}_{H \sim Q} \mathbb{E}_{\gamma \sim G} [\mathcal{R}^{e(\gamma)}(h)] \leq \mathbb{E}_{H \sim Q} [\hat{\mathcal{R}}_{\Pi_\phi}(h)] + \sqrt{\frac{KL(Q||M) + \ln(1/\delta)}{2n}} + \sqrt{\frac{\rho}{2}}. \quad (11)$$



**Interpretation.** Minimizing the empirical *policy-averaged* risk controls the risk under *any* test environment mixture  $G$  in a KL-ball of radius  $\rho$  around the learned policy  $\Pi_\phi$ , independently of which added robust penalty (IRM/REx/GroupDRO/CORAL) during training.

### 4.3 Entropic risk duality (GroupDRO view)

**Proposition 4.5** (KL-ball DRO equals entropic risk). *Fix  $h$  and write  $r_\gamma = \mathcal{R}^{e(\gamma)}(h)$ . For any  $\rho > 0$ ,*

$$\sup_{G: KL(G\|\Pi_\phi) \leq \rho} \mathbb{E}_G[r_\gamma] = \inf_{\lambda > 0} \frac{\rho + \log \mathbb{E}_{\gamma \sim \Pi_\phi} \exp\{\lambda r_\gamma\}}{\lambda}.$$

Moreover, if  $\Pi_\phi$  has finite support of size  $k$  then, for any  $\lambda > 0$ ,

$$\max_\gamma r_\gamma \leq \frac{1}{\lambda} \log\left(\frac{1}{k} \sum_\gamma e^{\lambda r_\gamma}\right) + \frac{\log k}{\lambda}.$$

Hence *max* risk is upper-bounded by a log-sum-exp (entropic) risk—**not** by “mean +  $\sqrt{\text{Var}}$ ” in general.

### 4.4 Loss-vs-likelihood averaging (safe replacement)

The optimization objective of our UAED framework is the minimization of the policy-averaged risk,  $\mathcal{R}_{\Pi_\phi}(h)$ . This risk corresponds to the Loss Averaging (LoA) objective. In this section, we analyze the relationship between this quantity and Log-Likelihood Averaging (LA). This comparison demonstrates that the LA quantity serves as a safe, variance controlled proxy for the LoA objective.

**Proposition 4.6** (LoA vs. LA gap is variance-controlled). *Let  $\ell_\gamma = \mathbb{E}_{z'|z} \ell(h, z')$  for fixed  $(h, z)$  and denote  $\mu = \mathbb{E}_{\Pi_\phi} \ell_\gamma$ . We define:*

$$\mathcal{L}_{\text{LoA}} = \mathbb{E}_{\Pi_\phi} \ell_\gamma, \quad \mathcal{L}_{\text{LA}} = -\log \mathbb{E}_{\Pi_\phi} e^{-\ell_\gamma}.$$

If  $\ell_\gamma - \mu$  is sub-Gaussian with proxy  $\sigma^2$ , then

$$0 \leq \mathcal{L}_{\text{LoA}} - \mathcal{L}_{\text{LA}} = \log \mathbb{E}_{\Pi_\phi} e^{-(\ell_\gamma - \mu)} \leq \sigma^2/2.$$

**Interpretation.** The proposition demonstrates a key theoretical advantage: minimizing  $\mathcal{L}_{\text{LA}}$  is a safe replacement for minimizing  $\mathcal{L}_{\text{LoA}}$ , as the two are equivalent up to a term proportional to the loss variance  $\sigma^2$ . This bounds mathematically confirms that optimizing an objective that is close to  $\mathcal{L}_{\text{LoA}}$  (the policy-averaged risk) implicitly introduces a regularization against cross-environment loss variance. We will further provide empirical evidence in the experiments.

While the preceding theoretical analyses rely on the loss  $l$  being bounded or sub-Gaussian, we can generalize these results to the sub-gamma condition covering a broader class of potentially heavy-tailed losses.

*Remark 4.7* (Sub-gamma variant). If  $\ell$  is sub-gamma with variance proxy  $\sigma^2$  and scale  $c$ , then in Theorem 4.4 the PAC-Bayes term becomes its standard sub-gamma analogue, and the robustness term in Lemma 4.3 is replaced by  $\inf_{\lambda \in (0, 1/c)} \{\rho/\lambda + \psi(\lambda)\}$  with  $\psi(\lambda) \leq \frac{\sigma^2 \lambda^2}{2(1-c\lambda)}$ .

### 4.5 Optimization Strategy

Our UAED framework requires minimizing a joint objective function that concurrently optimizes the model parameters  $\theta$  (for hypothesis  $h_\theta$  and the policy parameters  $\phi$  (for environment distribution  $\Pi_\phi$ ), as follows:

$$L(\theta, \phi) = \mathbb{E}_{\gamma \sim \Pi_\phi} \mathbb{E}_{z \sim P} \ell(h_\theta, T_\gamma(z)) + \eta \mathcal{P}_{\text{robust}}(\theta; \Pi_\phi) + \beta \text{KL}(\Pi_\phi \| \Pi_0), \quad (12)$$

where  $\mathcal{P}_{\text{robust}}$  represents the robust penalty (e.g., the IRM penalty). Due to the coupled nature of  $\theta$  and  $\phi$  within the objective, we adopt an alternating optimization scheme. At each training step, we sample a mini-batch and update the policy parameters  $\phi$  while holding  $\theta$  fixed, and then update  $\theta$  while keeping  $\phi$  fixed. Section A.7 provides a convergence analysis of our optimization strategy.

## 5 Experiments

We evaluate our universal adaptive framework across both synthetic benchmarks and real-world datasets, demonstrating consistent improvements across diverse spurious correlation types.

**Datasets.** We use three standard benchmarks. COLORED-MNIST (Arjovsky et al., 2019) is a binary digit task ( $< 5$  vs.  $\geq 5$ ) with spurious color correlations controlled by  $e \in [0, 1]$ ; models train on  $e \in \{0.1, 0.2\}$  and test on  $e = 0.9$ . ROTATED-MNIST (Ghifary et al., 2015) uses the same task with rotations  $\{0^\circ, 45^\circ, 90^\circ, 135^\circ, 180^\circ\}$ , training on  $\{0^\circ, 90^\circ\}$ . WaterBirds (Wah et al., 2011; Sagawa et al., 2020b) is constructed by compositing CUB-200-2011 foreground birds onto background scenes to induce spurious correlations between label and background; it defines four groups by (bird type, background).

**Implementation Details.** For the synthetic benchmarks, we use 3-layer MLPs with 256 hidden units, while for WaterBirds we employ ResNet18 and ResNet50 pretrained on IMAGENET (Krizhevsky et al., 2012). Models are optimized with an Adam optimizer (Kingma and Ba, 2015) with a learning rate of  $10^{-4}$ . The adaptive policy applies a learning rate multiplier of  $100\times$  for continuous and  $10\times$  for discrete settings. We sample environments using 5 Monte Carlo samples for synthetic datasets and 3 for WaterBirds, and report worst-case accuracy across environments or groups for evaluation.

### 5.1 A-IRM Results on Synthetic Benchmarks

**A-IRM vs. standard IRM.** Table 1 demonstrates that our A-IRM significantly outperforms standard IRM on both synthetic benchmarks, with remarkable improvements on ROTATED-MNIST (+28.4%). The lower variance of A-IRM indicates more stable training.

Table 1: A-IRM results on COLORED-MNIST and ROTATED-MNIST.

| Dataset | Method                       | Worst-Case Acc                       | Gain   |
|---------|------------------------------|--------------------------------------|--------|
| C-MNIST | IRM ( $e \in \{0.1, 0.2\}$ ) | $66.8 \pm 2.9$ (%)                   | —      |
|         | A-IRM                        | <b><math>72.3 \pm 1.6</math></b> (%) | +5.5%  |
| R-MNIST | IRM ( $\{0, 90\}$ )          | $65.8 \pm 0.7$ (%)                   | —      |
|         | A-IRM                        | <b><math>94.2 \pm 0.2</math></b> (%) | +28.4% |

**Comparison with IRM variants.** Moreover, A-IRM achieves competitive performance with specialized IRM variants and matches the performance of manually optimized oracle IRM setup (see Table 2). Specifically, A-IRM achieves the highest accuracy on the challenging COLORED-MNIST test set ( $e = 0.9$ ). Its performance matches the manually-tuned IRM ( $e \in \{0.2, 0.8\}$ ), demonstrating that the adaptive policy automatically discovers optimal environments.

Table 2: Comparison with IRM variants on COLORED-MNIST (test  $e = 0.9$ ). The gains are compared to IRM.

| Method                       | Test Acc (%)                     | Gain         |
|------------------------------|----------------------------------|--------------|
| ERM                          | $17.1 \pm 0.6$                   | -74.4%       |
| IRM (Arjovsky et al., 2019)  | $66.8 \pm 2.9$                   | —            |
| Meta-IRM (Bae et al., 2021)  | $70.5 \pm 0.9$                   | +5.5%        |
| BIRM (Lin et al., 2022)      | $69.8 \pm 1.2$                   | +4.5%        |
| IRM ( $e \in \{0.2, 0.8\}$ ) | $72.2 \pm 0.5$                   | +8.1%        |
| A-IRM (Ours)                 | <b><math>72.3 \pm 1.6</math></b> | <b>+8.2%</b> |

### 5.2 Environment Discovery Analysis

Next, we analyze the environment policies learned by A-IRM on the synthetic benchmarks to gain insight into its superior performance. The key finding is that A-IRM successfully identifies the most informative environment distributions for robust learning, which often differ substantially from those used in standard, fixed-IRM settings. Fig. 2 and Fig. 3 show the analyses on COLORED-MNIST and ROTATED-MNIST, respectively. We observe that, on COLORED-MNIST, A-IRM discovers that intermediate correlations ( $e \approx 0.35$ ) provide the optimal learning signal—avoiding both uninformative similar environments and conflicting diverse environments. Meanwhile, on ROTATED-MNIST, A-IRM learns approximately uniform distribution over all rotations, leading to robust performance across all test angles while IRM catastrophically fails on unseen rotations.

### 5.3 Theoretical Validation: Variance Regularization

Our theoretical analysis (Theorem 4.6) predicts that A-IRM’s loss-averaging objective implicitly minimizes cross-environment variance. Fig. 4(left) provides striking empirical validation: A-IRM achieves approximately  $100\times$  lower final loss variance ( $\sim 10^{-4}$ ) compared to IRM ( $\sim 10^{-2}$ ). This dramatic variance reduction explains why A-IRM discovers solutions that perform consistently across all environments rather than overfitting to specific ones.

### 5.4 Results on Waterbirds Dataset

Table 3: Performance on waterBirds dataset. All adaptive variants improve upon their baselines.

| Method     | Type     | ResNet18    |             | ResNet50    |             |
|------------|----------|-------------|-------------|-------------|-------------|
|            |          | Mean        | Worst       | Mean        | Worst       |
| ERM        | Baseline | 85.2        | 60.0        | 87.1        | 63.2        |
| GroupDRO   | Baseline | 88.0        | 78.0        | 93.2*       | 86.0*       |
| A-GroupDRO | Adaptive | <b>88.8</b> | <b>78.0</b> | <b>92.1</b> | <b>87.8</b> |
| REx/VREx   | Baseline | 87.0        | 74.0        | —           | —           |
| A-REx      | Adaptive | <b>88.8</b> | <b>78.0</b> | <b>92.7</b> | 80.5        |
| CORAL      | Baseline | 86.0        | 73.0        | —           | —           |
| A-CORAL    | Adaptive | <b>86.0</b> | <b>75.0</b> | <b>90.8</b> | 82.9        |

\*GroupDRO baseline from Sagawa et al. (2020a).

Results on WaterBirds dataset (Table 3) demonstrate that our UAED framework provides a consistent performance boost across multiple robust learning methods. In all tested scenarios (ERM, GroupDRO, REx, CORAL) and architectures (ResNet18, ResNet50), the Adaptive (A-) variants match or significantly exceed the worst-group accuracy of their vanilla baselines. Notably, A-GroupDRO achieves the highest worst-case accuracy of 87.8% (ResNet50), surpassing the strong GroupDRO baseline (86.0%) and confirming the efficacy of adaptively learning environment distributions—even when known groups exist. This universal improvement shows that UAED functions as a generic enhancement for targeting and improving generalization to the most vulnerable data subgroups. For additional experiments on PACS dataset (Li et al., 2017), see Section C.

### 5.5 Comparison with State-of-the-Art

To contextualize the performance of our UAED framework, Table 4 compares the worst-group accuracy of our top performing variant, A-GroupDRO, against state-of-the-art (SOTA) methods on the WaterBirds dataset using ResNet50. Our A-GroupDRO variant achieves a highly competitive worst-group accuracy of 87.8%, positioning the UAED framework favorably

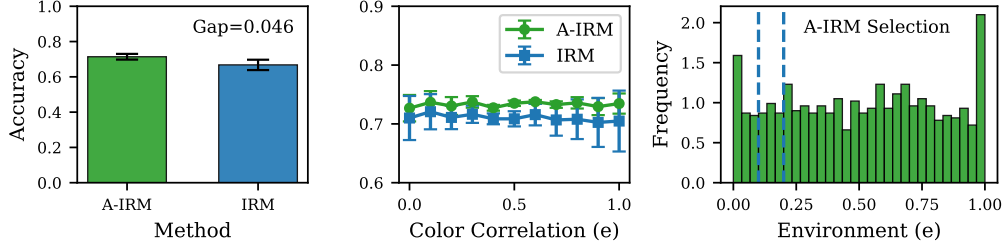


Figure 2: A-IRM environment discovery on COLORED-MNIST. **(Left)** Worst-case test accuracy—A-IRM achieves 72.3% vs IRM’s 66.8%, a 5.5% improvement. **(Middle)** Test accuracy across different color correlations shows A-IRM’s flatter profile, indicating successful invariant learning. **(Right)** A-IRM discovers intermediate correlations ( $e \approx 0.35$ ) rather than the extremes used by standard IRM (blue dashed lines at 0.1, 0.2).

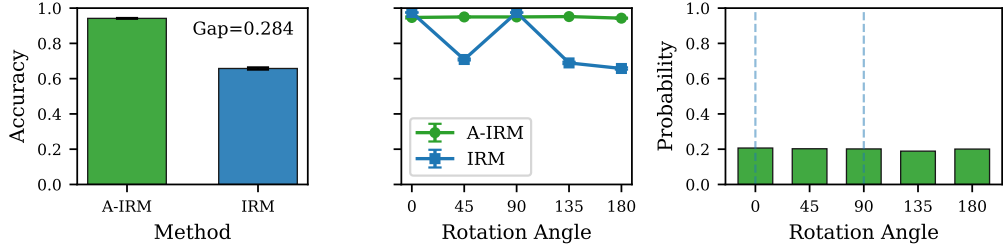


Figure 3: A-IRM environment discovery on ROTATED-MNIST. **(Left)** Worst-case accuracy comparison shows 28.4% improvement (94.2% vs 65.8%). **(Middle)** Test accuracy across rotation angles—IRM fails catastrophically on unseen rotations (45°, 135°, 180°) while A-IRM maintains consistent performance. **(Right)** A-IRM learns approximately uniform distribution over all rotations, unlike IRM which only uses 0° and 90° (indicated by dashed lines).

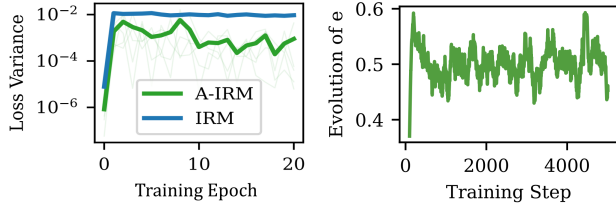


Figure 4: Evolution of the loss variance of A-IRM and IRM **(Left)**, and the environment  $e$  of A-IRM **(Right)** on COLORED-MNIST.

Table 4: Comparison with SOTA methods on Waterbirds (worst-group accuracy) using ResNet50. <sup>†</sup>Uses additional concept discovery. <sup>‡</sup>Requires retraining last layer.

| Method                          | Worst-Group Acc     |
|---------------------------------|---------------------|
| GroupDRO (Sagawa et al., 2020a) | 86.0 %              |
| JTT (Liu et al., 2021a)         | 86.7 %              |
| DISC (Li et al., 2023)          | 88.0 % <sup>†</sup> |
| DFR (Kirichenko et al., 2023)   | 91.2% <sup>‡</sup>  |
| A-GroupDRO ( <b>Ours</b> )      | 87.8 %              |

among specialized SOTA methods. Crucially, while methods like GroupDRO, JTT, DISC, and DFR either rely on external group annotations, require complex two-stage training procedures, or introduce architecture modifications, A-GroupDRO achieves result using a simple, single-stage, alternating optimization over parameterized transformations. This competitive performance, achieved without the need of pre-defined groups or significant procedure overhead, highlights the advantage of environment adaptivity as a generic,

assumption-free enhancement to robust learning.

Table 5: Ablation study on A-IRM (COLORED-MNIST).

| Method                  | Worst-Case Acc (%) | Mean $e$ Discovered |
|-------------------------|--------------------|---------------------|
| A-IRM (full)            | $72.3 \pm 1.6$     | $0.354 \pm 0.023$   |
| w/o hierarchical policy | $68.9 \pm 2.1$     | $0.287 \pm 0.041$   |
| w/o KL regularization   | $67.5 \pm 2.8$     | $0.198 \pm 0.018$   |
| w/ fixed variance       | $69.5 \pm 1.9$     | $0.312 \pm 0.015$   |
| Different MC samples:   |                    |                     |
| K=1                     | $70.1 \pm 2.0$     | $0.341 \pm 0.029$   |
| K=5 (default)           | $71.4 \pm 1.6$     | $0.354 \pm 0.023$   |
| K=10                    | $71.6 \pm 1.5$     | $0.356 \pm 0.022$   |

## 5.6 Ablation Studies

Ablation studies across COLORED-MNIST (Table 5) and WaterBirds (Table 7) confirm the necessity of the core components of our UAED framework. In particular, the hierarchical policy structure is crucial, as its removal leads to a 3.4% drop in A-IRM’s worst-case accuracy, underscoring its role in managing the large environment space. Furthermore, the KL regularization proves essential, preventing the environment policy  $\Pi_\phi$  from collapsing onto overly simple or challenging single environments, evidenced by a significant drop (4.8% for A-IRM) and a lower mean discovered environment factor ( $0.354 \rightarrow 0.198$ ). In addition, the number of MC samples,  $K$ , shows diminishing returns, with  $K = 5$  striking an optimal balance. Finally, the analysis of discovered environments (Table 5) reveals that UAED’s

Table 6: Discovered environment statistics across methods and datasets

| Method     | Dataset       | Discovered Pattern   |
|------------|---------------|--|
| A-IRM      | COLORED-MNIST | Concentrates on $e \approx 0.35$ : intermediate correlations that balance breaking spurious patterns with maintaining learnability |
| A-IRM      | ROTATED-MNIST | Uniform distribution over all rotations: maximizes diversity for invariance  |
| A-GroupDRO | WaterBirds    | Focuses on minority groups with high spurious correlation  |
| A-REx      | WaterBirds    | Discovers maximally diverse environments while maintaining stability   |
| A-CORAL    | WaterBirds    | Identifies environments with distinct feature distributions  |

Table 7: Ablation study on adaptive methods for WaterBirds (ResNet18).

| Method     | Configuration         | Mean Acc (%) | Worst-group Acc (%) |
|------------|-----------------------|--------------|---------------------|
| A-GroupDRO | Full model            | 88.8         | 78.0                |
|            | w/o adaptive policy   | 88.0         | 78.0                |
|            | w/o KL regularization | 87.5         | 75.2                |
| A-REx      | Full model            | 88.8         | 78.0                |
|            | w/o variance penalty  | 87.2         | 74.5                |
|            | Fixed environments    | 87.0         | 74.0                |
| A-CORAL    | Full model            | 86.0         | 75.0                |
|            | w/o MMD alignment     | 85.5         | 73.2                |
|            | Fixed environments    | 86.0         | 73.0                |

policy is not random but interpretable and objective-specific: A-IRM on COLORED-MNIST concentrates on an intermediate correlation ( $e \approx 0.35$ ) that maximizes the IRM penalty, while A-GroupDRO on WaterBirds focuses its probability mass on the known minor groups, effectively reproducing the ideal GroupDRO behavior without explicit group labels. These results validate our design choices and further confirm that the adaptive policy successfully discovers robust and meaningful environmental variations.

### 5.7 Detailed Environment Discovery Analysis

The detailed analysis of the learned environment policies (Table 6) provides crucial, interpretable evidence for the efficacy of our UAED framework. The resulting environment distributions are not generic but are precisely tailored to the objective of the base robust method. This validation confirms that our framework successfully finds the ideal challenging distribution  $\Pi_\phi$  that maximizes the specific robustness criteria for each baseline. This is evidenced by the policy’s ability to find the maximal diversity required for general invariance tasks (like ROTATED-MNIST) or to learn to focus on the most critical data subgroups (like WaterBirds), effectively automating complex environment specification without relying on external annotations.

### 5.8 Discussion: Why Universal Adaptive Discovery Works

Our experiments highlight three factors behind the effectiveness of adaptive environment discovery. (1)

*Method–Environment Alignment:* Different robust objectives benefit from different environments—e.g., A-GroupDRO discovers worst-case scenarios, while A-REx favors diverse yet learnable ones. The adaptive framework lets each method find the environments best suited to its objective. (2) *Continuous Exploration:* Unlike fixed specifications, adaptive discovery updates the environment distribution in response to model weaknesses, creating a curriculum that progressively challenges the learner. (3) *Implicit Regularization:* The hierarchical Bayesian prior and KL regularization discourage overfitting to particular environments while promoting exploration of diverse conditions.

## 6 Conclusion

We introduced *Universal Adaptive Environment Discovery* (UAED), a unified framework that shifts robust learning from relying on fixed environments to discovering them adaptively. Our work shows that diverse robust methods such as IRM, REX, GroupDRO, and CORAL can all benefit from adaptive environment discovery, supported by theoretical guarantees of distributionally robust generalization and empirical validation on synthetic and real-world data. Beyond empirical and theoretical contributions, we provide the conceptual insight that environment specification and robust objectives are complementary aspects of the same problem. These results suggest that the future of robust learning lies not in designing new objectives or manually defining environments, but in methods that jointly discover both—making environment discovery adaptive and automatic to provide a principled path toward truly robust machine learning systems.

**Limitations and future work.** Adaptive methods incur higher computational cost by sampling  $K$  environments per batch, though this overhead can be mitigated via parallelization. While environments are learned, the transformation family  $T_\gamma$  must still be specified, suggesting future work on learning transformation functions directly. Finally, our theory guarantees robustness within KL-balls of the learned distribution, but connecting these guarantees to worst-case out-of-distribution (OOD) performance remains an open question.



## References

- K. Ahuja, K. Shanmugam, K. Varshney, and A. Dhurandhar. Invariant Risk Minimization Games. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 145–155. PMLR, 13–18 Jul 2020.
- M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz. Invariant Risk Minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- J.-H. Bae, I. Choi, and M. Lee. Meta-learned Invariant Risk Minimization. *arXiv preprint arXiv:2103.12947*, 2021.
- Y. Bengio, J. Louradour, R. Collobert, and J. Weston. Curriculum Learning. In *Proceedings of the 26th International Conference on Machine Learning*, pages 41–48, 2009.
- G. Benton, M. Finzi, P. Izmailov, and A. G. Wilson. Learning Invariances in Neural Networks from Training Data. In *Advances in Neural Information Processing Systems*, volume 33, pages 17605–17616, 2020.
- P. G. Bissiri, C. C. Holmes, and S. G. Walker. A general framework for updating belief distributions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(5):1103–1130, Feb. 2016.
- E. Creager, J.-H. Jacobsen, and R. Zemel. Environment Inference for Invariant Learning. In *International Conference on Machine Learning*, pages 2189–2200. PMLR, 2021.
- E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le. AutoAugment: Learning Augmentation Strategies From Data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 113–123, 2019.
- E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le. RandAugment: Practical Automated Data Augmentation with a Reduced Search Space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 702–703, 2020.
- Z. Deng, F. Ding, C. Dwork, R. Hong, G. Parmigiani, P. Patil, and P. Sur. Representation via Representations: Domain Generalization via Adversarially Learned Invariant Representations. *arXiv preprint arXiv:2006.11478*, 2020.
- C. Finn, P. Abbeel, and S. Levine. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. In *International Conference on Machine Learning*, pages 1126–1135. PMLR, 2017.
- Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-Adversarial Training of Neural Networks. In *Journal of Machine Learning Research*, volume 17, pages 1–35, 2016.
- I. Gao, S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. Out-of-Domain Robustness via Targeted Augmentations. In *International Conference on Machine Learning*, pages 10799–10820. PMLR, 2023.
- A. Gelman, J. B. Carlin, H. S. Stern, and D. B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, 1995.
- P. Germain, A. Habrard, F. Laviolette, and E. Morvant. A PAC-Bayesian Approach for Domain Adaptation with Specialization to Linear Classifiers. In *International Conference on Machine Learning*, pages 738–746. PMLR, 2013.
- M. Ghifary, W. B. Kleijn, M. Zhang, and D. Balduzzi. Domain Generalization for Object Recognition With Multi-Task Autoencoders. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A Kernel Two-Sample Test. *Journal of Machine Learning Research*, 13(25):723–773, 2012.
- I. Gulrajani and D. Lopez-Paz. In Search of Lost Domain Generalization. In *International Conference on Learning Representations*, 2021.
- I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. beta-VAE: Learning basic visual concepts with a constrained variational framework. In *International Conference on Learning Representations*, 2017.
- E. Jang, S. Gu, and B. Poole. Categorical Reparameterization with Gumbel-Softmax. In *International Conference on Learning Representations*, 2017.
- M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. *An introduction to variational methods for graphical models*, page 105–161. MIT Press, Cambridge, MA, USA, 1999. ISBN 0262600323.
- D. Kim, Y. Yoo, S. Park, J. Kim, and J. Lee. SelfReg: Self-Supervised Contrastive Regularization for Domain Generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9619–9628, October 2021.
- D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations*, 2015.
- P. Kirichenko, P. Izmailov, and A. G. Wilson. Last Layer Re-Training is Sufficient for Robustness to Spurious Correlations. In *The Eleventh International Conference on Learning Representations*, 2023.

- A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012.
- D. Krueger, E. Caballero, J.-H. Jacobsen, A. Zhang, J. Binas, D. Zhang, R. Le Priol, and A. Courville. Out-of-Distribution Generalization via Risk Extrapolation (REx). In *International Conference on Machine Learning*, pages 5815–5826. PMLR, 2021.
- D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales. Deeper, Broader and Artier Domain Generalization, 2017.
- Y. Li, H. Han, S. Shan, and X. Chen. DISC: Learning From Noisy Labels via Dynamic Instance-Specific Selection and Correction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 24070–24079, June 2023.
- Y. Lin, H. Dong, H. Wang, and T. Zhang. Bayesian invariant risk minimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16021–16030, 2022.
- E. Z. Liu, B. Haghighi, A. S. Chen, A. Raghunathan, P. W. Koh, S. Sagawa, P. Liang, and C. Finn. Just Train Twice: Improving Group Robustness without Training Group Information. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6781–6792. PMLR, 18–24 Jul 2021a.
- J. Liu, Z. Hu, P. Cui, B. Li, and Z. Shen. Heterogeneous Risk Minimization. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 6804–6814. PMLR, 18–24 Jul 2021b.
- M. Matymov, B.-H. Tran, M. Kampffmeyer, M. Heinonen, and M. Filippone. Optimizing data augmentation through bayesian model selection. *arXiv preprint arXiv:2505.21813*, 2025.
- D. A. McAllester. PAC-Bayesian Model Averaging. In *Proceedings of the Twelfth Annual Conference on Computational Learning Theory*, pages 164–170, 1999.
- A. Nichol, J. Achiam, and J. Schulman. On First-Order Meta-Learning Algorithms. *arXiv preprint arXiv:1803.02999*, 2018.
- J. Peters, P. Bühlmann, and N. Meinshausen. Causal inference using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B*, 2016.
- E. Rosenfeld, P. Ravikumar, and A. Risteski. The Risks of Invariant Risk Minimization. In *International Conference on Learning Representations*, 2021.
- S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. Distributionally Robust Neural Networks for Group Shifts: On the Importance of Regularization for Worst-Case Generalization. In *International Conference on Learning Representations*, 2020a.
- S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. S. Liang. An Investigation of Why Overparameterization Exacerbates Spurious Correlations. In *ICML Workshop on Uncertainty & Robustness in Deep Learning*, 2020b.
- B. Sun and K. Saenko. Deep CORAL: Correlation Alignment for Deep Domain Adaptation. In *European Conference on Computer Vision*, pages 443–450. Springer, 2016.
- C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. 2011.
- H. Wang, S. Ge, Z. Lipton, and E. P. Xing. Learning Robust Global Representations by Penalizing Local Predictive Power. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le. Self-Training With Noisy Student Improves ImageNet Classification. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- H. Yao, Y. Wang, S. Li, L. Zhang, W. Liang, J. Zou, and C. Finn. Improving Out-of-Distribution Robustness via Selective Augmentation. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 25407–25437. PMLR, 17–23 Jul 2022.
- H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz. mixup: Beyond Empirical Risk Minimization. In *International Conference on Learning Representations*, 2018.

---

# Universal Adaptive Environment Discovery: Supplementary Materials

---

## A Theoretical Proofs

### A.1 Preliminaries

We use the following standard facts.

**Lemma A.1** (Hoeffding’s lemma). *If  $X \in [a, b]$ , then  $\log \mathbb{E} e^{\lambda(X - \mathbb{E}X)} \leq \frac{\lambda^2(b-a)^2}{8}$  for all  $\lambda \in \mathbb{R}$ .*

**Lemma A.2** (Donsker–Varadhan (DV)). *For probability measures  $G, \Pi$  and any measurable  $g$ ,  $\mathbb{E}_G[g] \leq \frac{KL(G\|\Pi) + \log \mathbb{E}_\Pi e^{\lambda g}}{\lambda}$ ,  $\forall \lambda > 0$ .*

### A.2 Proof of Theorem 4.2

Consider the composite bounded loss  $\tilde{\ell}(h, \phi; z) = \mathbb{E}_{\gamma \sim \Pi_\phi} \mathbb{E}_{z'|z} \ell(h, z') \in [0, 1]$ . Apply the classical PAC–Bayes bound (e.g., for bounded losses) to the joint hypothesis  $H = (h, \phi)$  with prior  $M = P \times \Pi_0$  and posterior  $Q$ : with probability  $\geq 1 - \delta$ ,

$$\mathbb{E}_{H \sim Q} \mathbb{E}_{z \sim P} \tilde{\ell}(H; z) \leq \mathbb{E}_{H \sim Q} \frac{1}{n} \sum_{i=1}^n \tilde{\ell}(H; z_i) + \sqrt{\frac{KL(Q\|M) + \ln(1/\delta)}{2n}}.$$

Identifying the terms with  $\mathcal{R}_{\Pi_\phi}$  and  $\hat{\mathcal{R}}_{\Pi_\phi}$  completes the proof.  $\square$

### A.3 Proof of Lemma 4.3

Let  $f \in [0, 1]$ . By Lemma A.2, for any  $\lambda > 0$ ,  $\mathbb{E}_G f \leq \frac{KL(G\|\Pi_\phi) + \log \mathbb{E}_{\Pi_\phi} e^{\lambda f}}{\lambda}$ . Write  $f = (f - \mathbb{E}_{\Pi_\phi} f) + \mathbb{E}_{\Pi_\phi} f$  and apply Lemma A.1 with  $a = 0, b = 1$ :  $\log \mathbb{E}_{\Pi_\phi} e^{\lambda(f - \mathbb{E}_{\Pi_\phi} f)} \leq \lambda^2/8$ . Hence  $\mathbb{E}_G f \leq \mathbb{E}_{\Pi_\phi} f + \frac{KL(G\|\Pi_\phi)}{\lambda} + \frac{\lambda}{8}$ . Optimizing over  $\lambda > 0$  gives  $\lambda^* = \sqrt{8 KL(G\|\Pi_\phi)}$  and the value  $\mathbb{E}_{\Pi_\phi} f + \sqrt{\frac{1}{2} KL(G\|\Pi_\phi)}$ .  $\square$

### A.4 Proof of Theorem 4.4

Fix  $Q$ . Let  $f_H(\gamma) = \mathcal{R}^{e(\gamma)}(h)$  for  $H = (h, \phi)$ . By Lemma 4.3,

$$\mathbb{E}_{H \sim Q} \mathbb{E}_{\gamma \sim G} f_H(\gamma) \leq \mathbb{E}_{H \sim Q} \mathbb{E}_{\gamma \sim \Pi_\phi} f_H(\gamma) + \sqrt{\frac{1}{2} KL(G\|\Pi_\phi)}.$$

Apply Theorem 4.2 to the first term and use  $KL(G\|\Pi_\phi) \leq \rho$ .  $\square$

### A.5 Proof of Proposition 4.5

By DV (Lemma A.2),

$$\sup_{G: KL(G\|\Pi_\phi) \leq \rho} \mathbb{E}_G r_\gamma = \inf_{\lambda > 0} \sup_G \frac{KL(G\|\Pi_\phi) - \rho + \log \mathbb{E}_{\Pi_\phi} e^{\lambda r_\gamma}}{\lambda} = \inf_{\lambda > 0} \frac{\rho + \log \mathbb{E}_{\Pi_\phi} e^{\lambda r_\gamma}}{\lambda},$$

where we used the Fenchel dual of the indicator of the KL–ball. For the finite-support bound, note that for uniform  $\Upsilon$  on the support ( $|\text{supp}| = k$ ),

$$\max_\gamma r_\gamma \leq \frac{1}{\lambda} \log \sum_\gamma \Upsilon(\gamma) e^{\lambda r_\gamma} + \frac{KL(\Upsilon\|\text{Unif})}{\lambda} = \frac{1}{\lambda} \log \left( \frac{1}{k} \sum_\gamma e^{\lambda r_\gamma} \right) + \frac{\log k}{\lambda}.$$

$\square$

## A.6 Proof of Proposition 4.6

By Jensen,  $\mathcal{L}_{\text{LA}} \leq \mathcal{L}_{\text{LoA}}$ , hence the gap is nonnegative. Let  $Y = \ell_\gamma - \mu$ ; then  $\mathcal{L}_{\text{LoA}} - \mathcal{L}_{\text{LA}} = \log \mathbb{E} e^{-Y}$ . If  $Y$  is sub-Gaussian with proxy  $\sigma^2$ , then  $\log \mathbb{E} e^{tY} \leq \sigma^2 t^2/2$  for all  $t \in \mathbb{R}$ . With  $t = -1$  this gives the stated upper bound.  $\square$

## A.7 Convergence Analysis of the Alternating Optimization Strategy in Section 4.5

Our joint objective function of the model parameters  $\theta$  (for hypothesis  $h_\theta$  and the policy parameters  $\phi$  (for environment distribution  $\Pi_\phi$ ) is defined follows:

$$L(\theta, \phi) = \mathbb{E}_{\gamma \sim \Pi_\phi} \mathbb{E}_{z \sim P} \ell(h_\theta, T_\gamma(z)) + \mathcal{P}(h_\theta, \Phi_\phi) + \beta \cdot \text{KL}(\Pi_\phi \| \Pi_0).$$

**Assumption A.3** (Smoothness and gradient noise).  $L$  is lower bounded and has  $L$ -Lipschitz gradient; the reparameterized-gradient estimators for  $(\theta, \phi)$  are unbiased with bounded variance; step sizes satisfy  $\sum_t \eta_t = \infty$ ,  $\sum_t \eta_t^2 < \infty$ .

**Theorem A.4** (Nonconvex alternating SGD). *Under Assumption A.3, alternating stochastic gradient updates on  $(\theta, \phi)$  satisfy*

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} [\|\nabla L(\theta_t, \phi_t)\|^2] = \mathcal{O}(T^{-1/2}).$$

*Proof.* Under Assumption A.3, the standard descent lemma for  $L$  with  $L$ -Lipschitz gradient yields, for the alternating update,

$$\mathbb{E}[L_{t+1}] \leq \mathbb{E}[L_t] - \frac{\eta_t}{2} \mathbb{E} \|\nabla L_t\|^2 + C\eta_t^2,$$

for some constant  $C$  depending on the gradient-noise variance. Summing and using  $\sum_t \eta_t = \infty$ ,  $\sum_t \eta_t^2 < \infty$  gives  $\frac{1}{\sum_{t=1}^T \eta_t} \sum_{t=1}^T \eta_t \mathbb{E} \|\nabla L_t\|^2 \leq \mathcal{O}(1/\sum_{t=1}^T \eta_t) + \mathcal{O}(\frac{\sum_{t=1}^T \eta_t^2}{\sum_{t=1}^T \eta_t})$ . With  $\eta_t = \eta/\sqrt{t}$  this becomes  $\mathcal{O}(T^{-1/2})$ .  $\square$

## B Extended Experimental Details

### B.1 Implementation Details

#### Architecture Details:

- ResNet18/ResNet50: Pretrained on IMAGENET, final layer replaced for binary classification
- Feature dimension: 512 for ResNet18, 2048 for ResNet50
- Dropout: 0.0 (following GroupDRO setup)

#### Training Details:

- Optimizer: Adam with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$
- Learning rate:  $10^{-4}$  for model,  $10^{-3}$  for policy parameters
- Weight decay:  $10^{-5}$
- Batch size: 128
- Number of epochs: 30
- Early stopping: Based on worst-group validation accuracy

#### Adaptive Policy Configuration:

- Continuous policy: Beta distribution with learnable  $\alpha, \beta$  parameters



- Policy network: 2-layer MLP with 64 hidden units
- KL regularization weight  $\alpha$ : 1.0
- Monte Carlo samples: 3 per batch
- Warm-up: 5 epochs before enabling adaptive policy

## B.2 Dataset Details

### Waterbirds:

- Training: 4,795 samples
- Validation: 1,199 samples
- Test: 5,794 samples
- Groups: 4 groups based on (bird type, background) combinations
- Group distribution: Highly imbalanced with smallest group having  $< 100$  training samples

**Assets & Licenses.** We use only publicly available research datasets and pretrained models:

- **CUB-200-2011** (Wah et al., 2011): license published by creators (non-commercial academic use).
- **Waterbirds** (Sagawa et al., 2020b): derived from CUB and Places; follows the respective research-use terms.
- **MNIST** (Colored/Rotated variants): generated from MNIST via our scripts; inherits MNIST’s research-use terms.
- **Pretrained ResNet-18/50** (PyTorch/torchvision): model weights and code under the PyTorch/torchvision license.

We release only *code/configs* to reproduce results (no new datasets or human/PII data). No additional consent was sought or required beyond the datasets’ published licenses/terms. The experiments contain no personally identifiable or offensive content.

## C Additional Experimental Results

### C.1 Results on PACS Dataset (Li et al., 2017)

Table 8: Results on PACS dataset

| Method     | Photo | Art   | Cartoon | Sketch | Average |
|------------|-------|-------|---------|--------|---------|
| A-GroupDRO | 85.84 | 77.09 | 97.66   | 80.40  | 85.25   |
| A-REx      | 85.45 | 74.06 | 97.49   | 79.10  | 84.16   |
| A-CORAL    | 84.80 | 73.95 | 97.11   | 79.05  | 83.73   |

## C.2 Additional Ablation Studies

Table 9: Effect of different transformation families on WaterBirds

| Transformation Type               | Mean Acc | Worst-Group Acc |
|-----------------------------------|----------|-----------------|
| Correlation strength (continuous) | 88.8     | 78.0            |
| Discrete groups                   | 87.9     | 76.5            |

## D Extended Related Work

### D.1 Connection to Meta-Learning

Our UAED framework shares conceptual similarities with the bi-level optimization common in meta-learning approaches (Finn et al., 2017; Nichol et al., 2018), as both involve learning a higher-level policy ( $\phi$ ) to guide a lower-level model ( $\theta$ ). However, their objectives differ: Meta-learning seeks adaptability—finding a strategy (e.g., an initialization) for fast learning on new tasks with few examples—while UAED seeks invariance by discovering the optimal distribution of synthetic environments necessary to enforce a robust predictor that generalizes across all potential distribution shifts.

### D.2 Connection to Curriculum Learning

Our adaptive framework UAED implicitly implements a form of curriculum learning (Bengio et al., 2009) by dynamically adjusting the environment policy  $\Pi_\phi$  throughout training. In the early training phase, the policy explores diverse environments, easing the initial learning task. During mid training, the focus shifts to more challenging yet informative and learnable environments, ensuring the model efficiently extracts the invariant signal. Finally, in the late training phase, the policy concentrates on the worst-case scenarios (e.g., maximizing risk variance for A-REx or the max risk for A-GroupDRO), effectively stress-testing the mature model to achieve maximal distributional robustness and pushing the invariant predictor to its generalization limits. This continuous, self-paced adjustment of the training distribution is key to the framework’s stability and superior performance.

### D.3 Domain Adaptation vs Environment Discovery

Our UAED framework differs fundamentally from traditional Domain Adaptation (DA) (Ganin et al., 2016). While DA relies on access to a specific target domain’s data to align features, UAED requires no target domain; instead, it actively discovers and generates multiple synthetic environments via parameterized kernels. Crucially, UAED optimizes for worst-case robustness across the learned environment distribution, aiming for general invariance, rather than targeting performance on a single, fixed distribution.