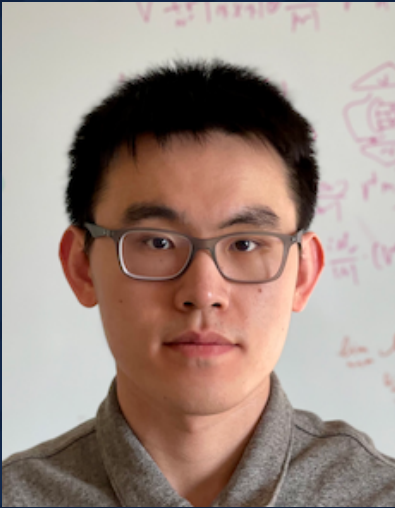




UNIVERSITY OF  
**ILLINOIS**  
URBANA-CHAMPAIGN

# Provable Domain Generalization via Invariant-Feature Subspace Recovery

ICML 2022



**Haoxiang Wang**  
PhD Candidate  
ECE, UIUC



**Haozhe Si**  
Master Student  
ECE, UIUC



**Bo Li**  
Assistant Professor  
Computer Science, UIUC





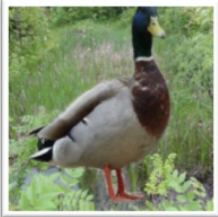

**Han Zhao**  
Assistant Professor  
Computer Science, UIUC

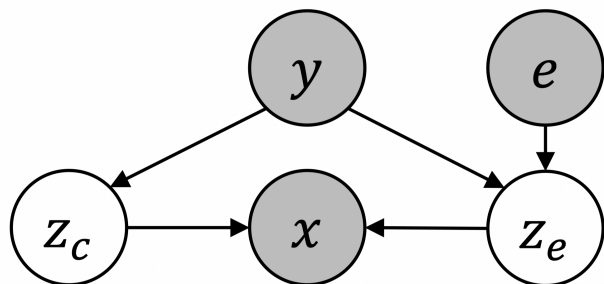
## Domain Generalization (OOD Generalization)

Perspective

## Spurious Correlation



		label: object	
		waterbird	landbird
spurious attribute: background	water background	 majority	 minority
	land background	 minority	 majority



In arbitrary training environment  $e = 1, 2, \dots, E$ , a sample  $(x, y, e)$  is generated by:

$$y = \begin{cases} 1, & \text{with probability } \eta \\ -1, & \text{otherwise} \end{cases}$$

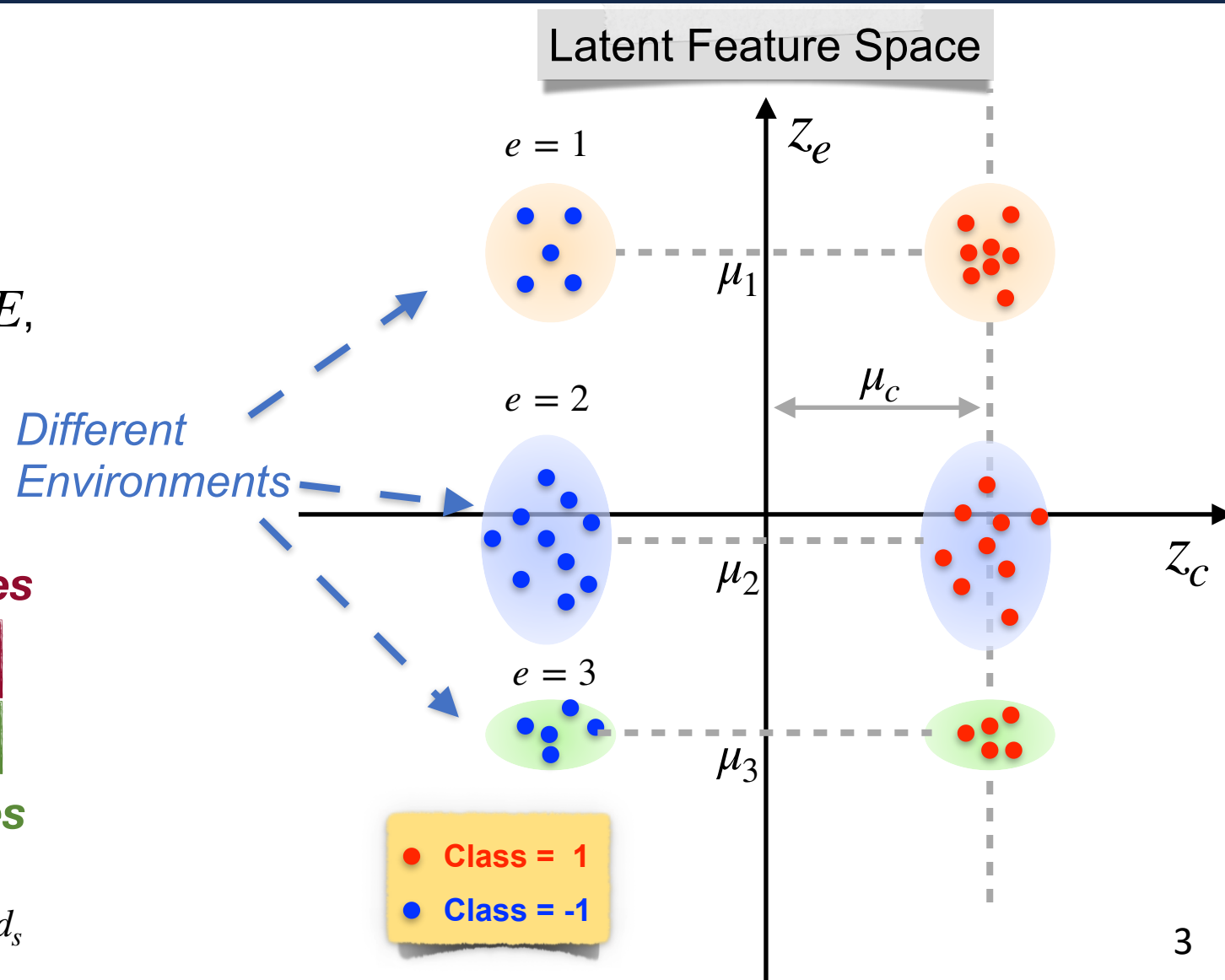
**Invariant Features**

$$z = \begin{bmatrix} z_c \\ z_e \end{bmatrix} \begin{cases} z_c \sim \mathcal{N}(y\mu_c, \sigma_c^2 I) \in \mathbb{R}^{d_c} \\ z_e \sim \mathcal{N}(y\mu_e, \sigma_e^2 I) \in \mathbb{R}^{d_s} \end{cases}$$

**Spurious Features**

$$x = Az_c + Bz_e \in \mathbb{R}^d,$$

$$\text{where, } d = d_c + d_s, A = \mathbb{R}^{d \times d_c}, B = \mathbb{R}^{d \times d_s}$$





IRM optimizes a bi-level objective over a feature extractor  $\Phi$  and a classifier  $\beta$  (assumed both to be **linear** here)

$$\min_{\Phi, \beta} \sum_{e \in [E]} \mathcal{R}^e(\Phi, \beta)$$

$$\text{s.t. } \beta \in \arg \min_{\beta} \mathcal{R}^e(\Phi, \beta) \quad \forall e \in [E]$$

**Goal of IRM**  $\rightarrow$  **Optimal Invariant Predictor (OIP)**

Given  $x = [A \ B] \begin{bmatrix} z_c \\ z_e \end{bmatrix}$ , an example of OIP is:

$$\Phi^*(x) = \begin{bmatrix} z_c \\ 0 \end{bmatrix} \text{ (keep invariant features only)}$$

$\beta^*$  is the optimal classifier w.r.t.  $\{(\Phi^*(x), y)\}$

**Assumption 1 [Mean].** For  $\{\mu_e\}_{e=1}^E$ , each element cannot be expressed as an affine combination of the rest.

**Assumption 2 [Covariance].** There exists a pair of distinct environments  $e, e' \in [E]$  s.t.  $\sigma_e \neq \sigma_{e'}$ .

**Linear Environment Complexity:** as  $E > d_s$ , the global optimum of IRM is guaranteed to be an optimal invariant predictor (*proved in [Rosenfeld et al. ICLR 2021]*).

**Convergence Issue:** IRM has no global convergence guarantee due to non-convexity.



*Goal of ISR: Find the feature subspace spanned by invariant features only.*

---

## Algorithm 1 ISR-Mean

---

**Input:** Data of all training environments,  $\{\mathcal{D}_e\}_{e \in [E]}$ .

**for**  $e = 1, 2, \dots, E$  **do**

    Estimate the sample mean of  $\{x | (x, y) \in \mathcal{D}_e, y = 1\}$   
    as  $\bar{x}_e \in \mathbb{R}^d$

**end for**

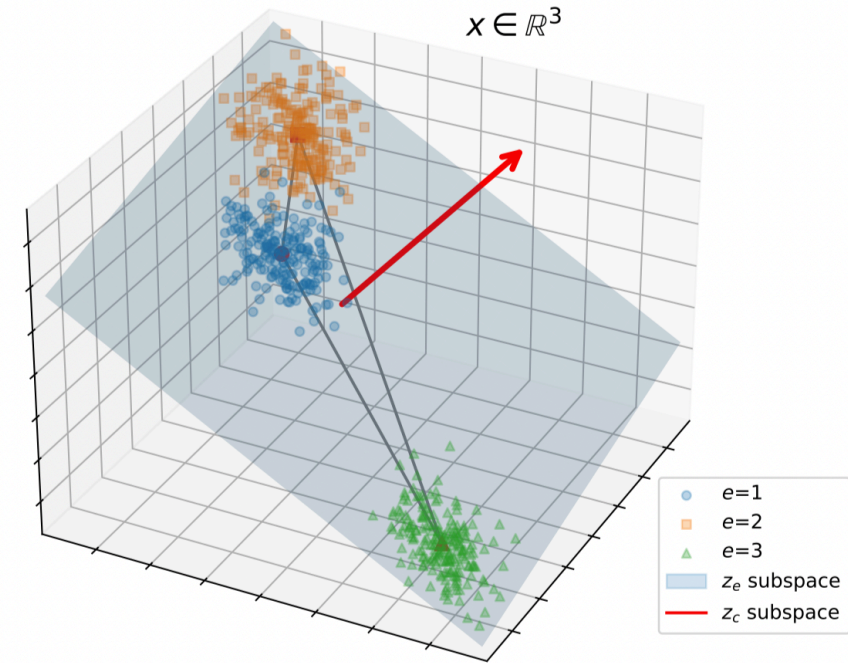
**I.** Construct a matrix  $\mathcal{M} \in \mathbb{R}^{E \times d}$  with the  $e$ -th row as  $\bar{x}_e^\top$   
for  $e \in [E]$

**II.** Apply PCA to  $\mathcal{M}$  to obtain eigenvectors  $\{P_1, \dots, P_d\}$   
with eigenvalues  $\{\lambda_1, \dots, \lambda_d\}$

**III.** Stack  $d_c$  eigenvectors with the lowest eigenvalues to  
obtain a transformation matrix  $P' \in \mathbb{R}^{d_c \times d}$

**IV.** Fit a linear classifier (with  $w \in \mathbb{R}^{d_c}$ ,  $b \in \mathbb{R}$ ) by ERM  
over all training data with transformation  $x \mapsto P'x$   
Obtain a predictor  $f(x) = w^\top P'x + b$

---



**Linear Environment Complexity:** Same as IRM.

**Global Convergence:** ISR-Mean enjoys global convergence guarantees.

**Less Assumptions than IRM:** Only Need Assumption 1 [Mean] & No need for Assumption 2 [Covariance].

---

## Algorithm 2 ISR-Cov

---

**Input:** Data of all training environments,  $\{\mathcal{D}_e\}_{e \in [E]}$ .

**for**  $e = 1, 2, \dots, E$  **do**

Estimate the sample covariance of  $\{x | (x, y) \in \mathcal{D}_e, y = 1\}$  as  $\Sigma_e \in \mathbb{R}^{d \times d}$

**end for**

**I.** Select a pair of environments  $e_1, e_2$  such that  $\Sigma_1 \neq \Sigma_2$ , and compute their difference,  $\Delta\Sigma := \Sigma_{e_1} - \Sigma_{e_2}$

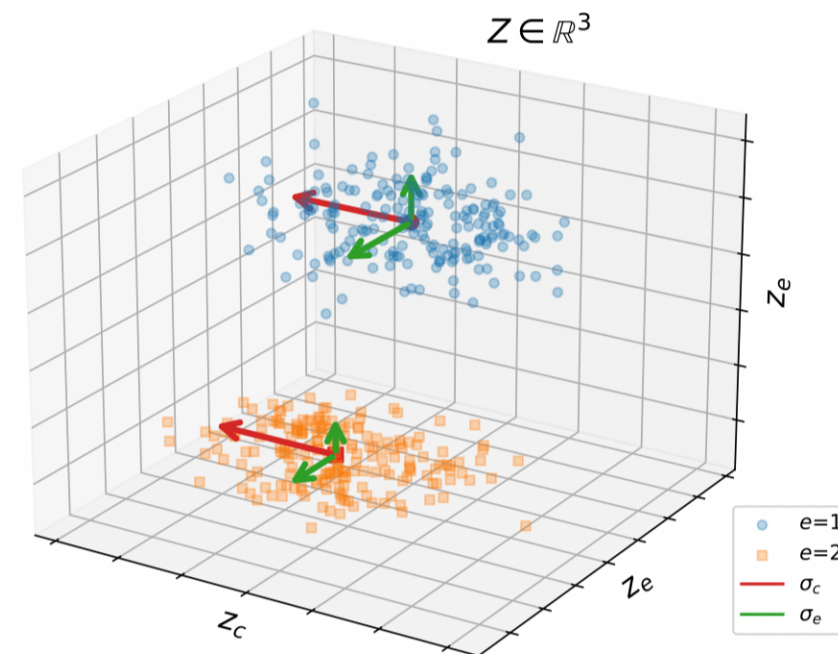
**II.** Eigen-decompose  $\Delta\Sigma$  to obtain eigenvectors  $\{P_1, \dots, P_d\}$  with eigenvalues  $\{\lambda_1, \dots, \lambda_d\}$

**III.** Stack  $d_c$  eigenvectors of eigenvalues with lowest absolute values to obtain a matrix  $P' \in \mathbb{R}^{d_c \times d}$

**IV.** Fit a linear classifier (with  $w \in \mathbb{R}^{d_c}$ ,  $b \in \mathbb{R}$ ) by ERM over all training data with transformation  $x \mapsto P'x$

Obtain a predictor  $f(x) = w^\top P'x + b$

---

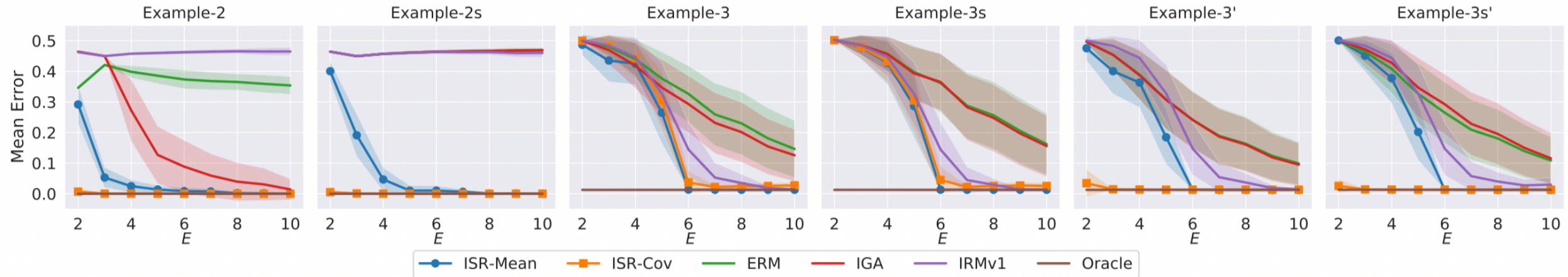


**$O(1)$  Environment Complexity:** Only needs 2 environments  $\longrightarrow$  optimal invariant predictor.

**Global Convergence:** ISR-Cov enjoys global convergence guarantees.

**Less Assumptions than IRM:** Only Need Assumption 2 [Covariance] & No need for Assumption 1 [Mean].






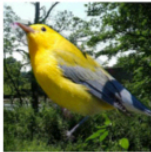




Test results on Linear Unit-Tests (first 4 plots) and its variants (last 2 plots), where  $d_c = 5$ ,  $d_s = 5$ , and  $E = 2, \dots, 10$ .  
 $P(Y|\mu_c)$  is invariant across environments.

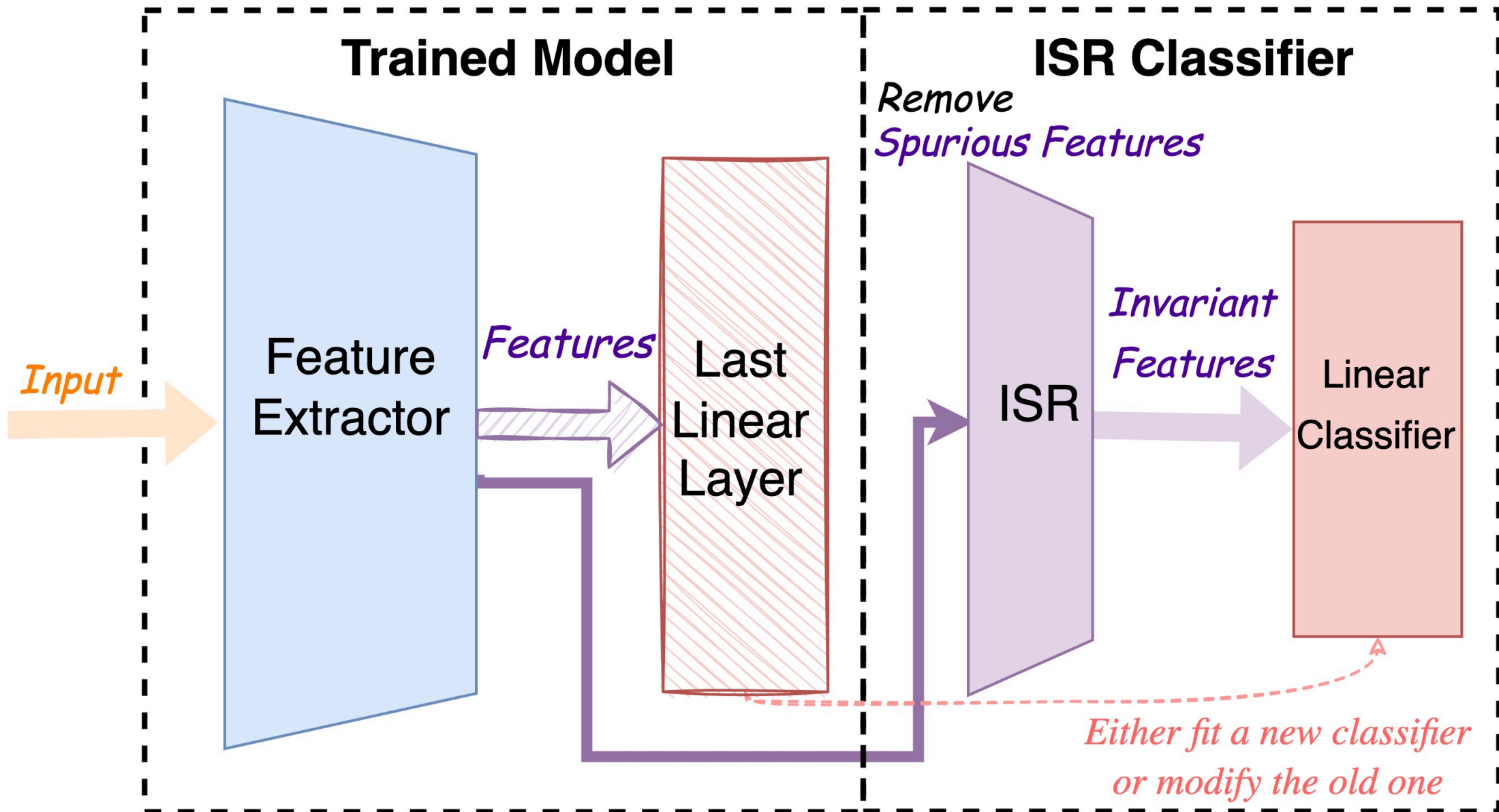
**ISR-Mean**      **Linear Environment Complexity:** Error is reduced to zero as  $E > d_s$   
**Better Performance than IRM:** Global convergence of ISR-Mean.

**ISR-Cov**      **O(1) Environment Complexity:** Error is reduced to zero as  $E \geq 2$  for datasets that satisfy Assumption 2 [Covariance].



**Benchmarks:** three datasets used by [Sagawa et al. ICLR 2020] to study the robustness of models against spurious correlations and group shifts.

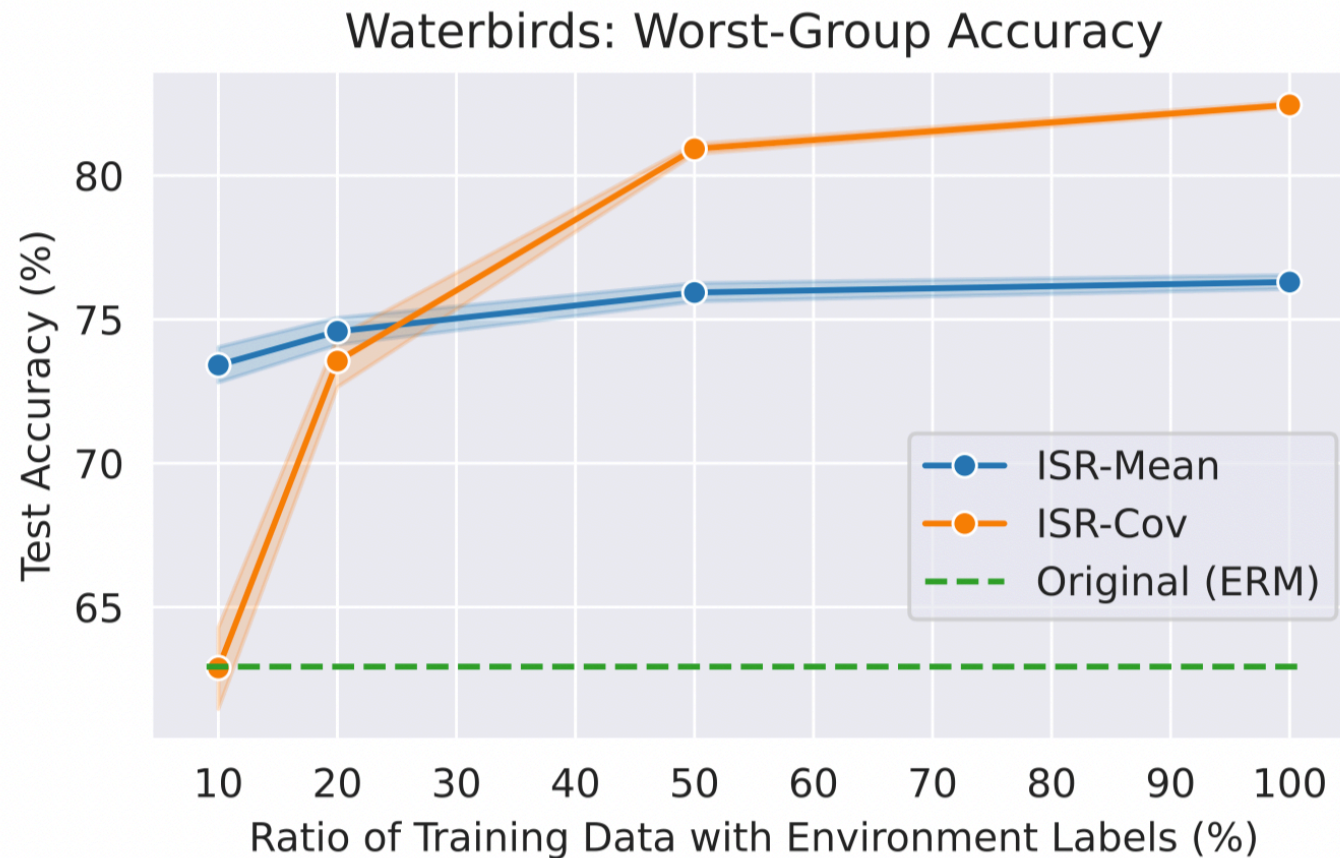
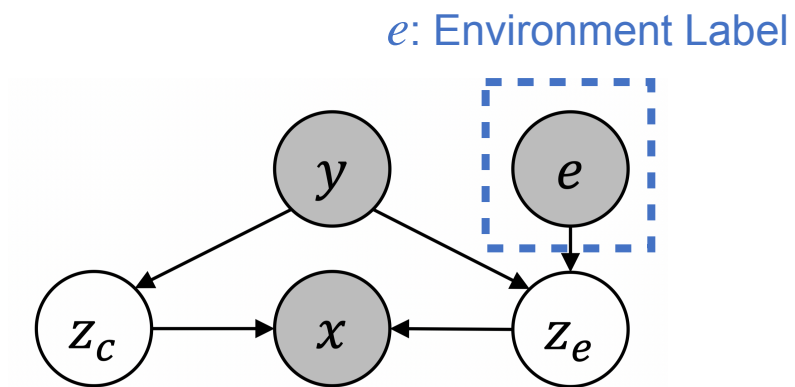
Common Training Examples				Test Examples	
Waterbird	Y: waterbird A: water background		Y: landbird A: land background		Y: landbird A: water background
					
CelebA	Y: blond hair A: female		Y: dark hair A: male		Y: blond hair A: male
					
MultiNLI	Y: contradiction A: has negation	(P) Abortive countryside revolts. (H) There is no revolt.	Y: entailment A: no negation	(P) The sacred is not mysterious to her. (H) The woman is familiar with the sacred.	Y: entailment A: has negation
					(P) Fixing current levels of damage would be impossible. (H) Fixing the damage could never be done.



Dataset	Backbone	Algorithm	Average Accuracy			Worst-Group Accuracy		
			Original	ISR-Mean	ISR-Cov	Original	ISR-Mean	ISR-Cov
Waterbirds	ResNet-50	ERM	86.66 $\pm$ 0.67	87.87 $\pm$ 0.80	<b>90.47<math>\pm</math>0.33</b>	62.93 $\pm$ 5.37	76.10 $\pm$ 1.11	<b>82.46<math>\pm</math>0.55</b>
		Reweighting	91.49 $\pm$ 0.46	<b>91.77<math>\pm</math>0.52</b>	91.63 $\pm$ 0.44	87.69 $\pm$ 0.53	88.02 $\pm$ 0.42	<b>88.67<math>\pm</math>0.55</b>
		GroupDRO	92.01 $\pm$ 0.33	91.74 $\pm$ 0.35	<b>92.25<math>\pm</math>0.27</b>	90.79 $\pm$ 0.47	90.42 $\pm$ 0.61	<b>91.00<math>\pm</math>0.45</b>
CelebA	ResNet-50	ERM	<b>95.12<math>\pm</math>0.34</b>	94.34 $\pm$ 0.12	90.12 $\pm$ 2.59	46.39 $\pm$ 2.42	55.39 $\pm$ 6.13	<b>79.73<math>\pm</math>5.00</b>
		Reweighting	<b>91.45<math>\pm</math>0.50</b>	91.38 $\pm$ 0.51	91.24 $\pm$ 0.35	84.44 $\pm$ 1.66	<b>90.08<math>\pm</math>0.50</b>	88.84 $\pm$ 0.57
		GroupDRO	<b>91.82<math>\pm</math>0.27</b>	91.82 $\pm$ 0.27	91.20 $\pm$ 0.23	88.22 $\pm$ 1.67	<b>90.95<math>\pm</math>0.32</b>	90.38 $\pm$ 0.42
MultiNLI	BERT	ERM	<b>82.48<math>\pm</math>0.40</b>	82.11 $\pm$ 0.18	81.28 $\pm$ 0.52	65.95 $\pm$ 1.65	72.60 $\pm$ 1.09	<b>74.21<math>\pm</math>2.55</b>
		Reweighting	<b>80.82<math>\pm</math>0.79</b>	80.53 $\pm$ 0.88	80.73 $\pm$ 0.90	64.73 $\pm$ 0.32	<b>67.87<math>\pm</math>0.21</b>	66.34 $\pm$ 2.46
		GroupDRO	<b>81.30<math>\pm</math>0.23</b>	81.21 $\pm$ 0.24	81.20 $\pm$ 0.24	78.43 $\pm$ 0.87	<b>78.95<math>\pm</math>0.95</b>	78.91 $\pm$ 0.75

- ISR classifiers can persistently **improve the worst-group accuracy** of trained models  
 → *ISR classifiers rely less on spurious features than original classifiers*
- The **average accuracy** of ISR classifiers is maintained around **the same level** as the original classifiers.





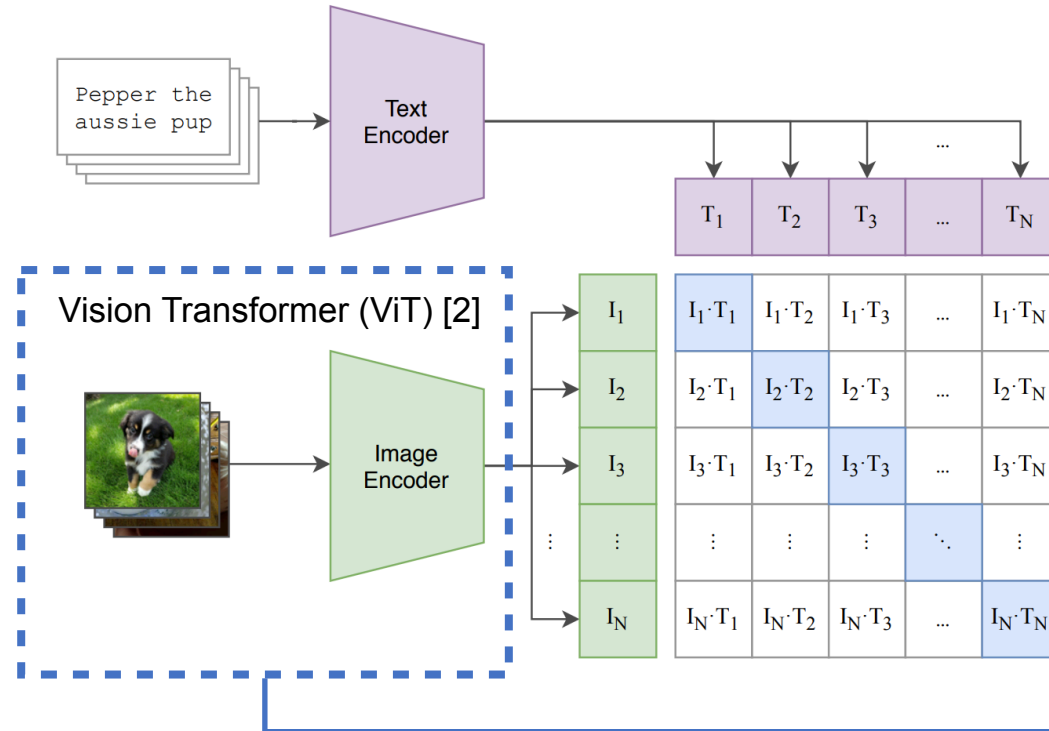
ISRs can be used in cases where only **a subset of** training samples have environment labels.

# ISRs for Pre-trained Feature Extractors (No need for neural net training!)

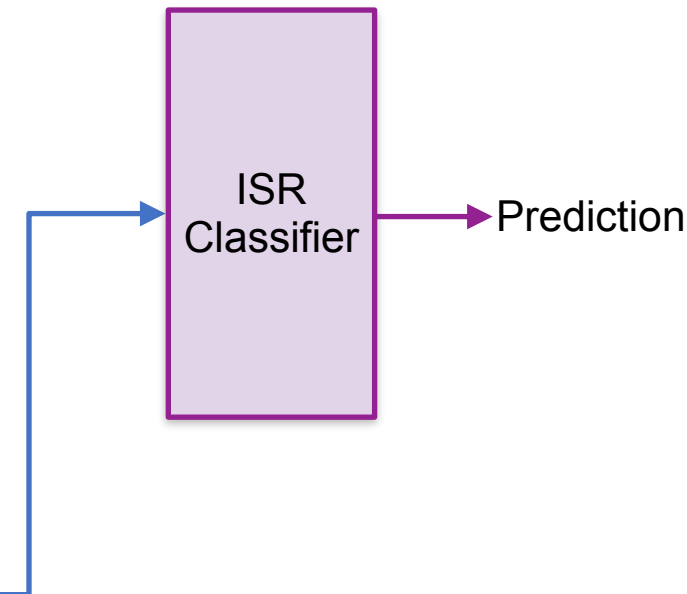


Dataset	Backbone	Algorithm	Average Accuracy			Worst-Group Accuracy		
			Linear Probing	ISR-Mean	ISR-Cov	Linear Probing	ISR-Mean	ISR-Cov
Waterbirds	CLIP (ViT-B/32)	ERM	$76.42 \pm 0.00$	<b><math>90.27 \pm 0.09</math></b>	$76.80 \pm 0.01$	$52.96 \pm 0.00$	<b><math>71.75 \pm 0.39</math></b>	$55.76 \pm 0.00$
		Reweighting	$87.38 \pm 0.09$	<b><math>88.23 \pm 0.12</math></b>	$88.07 \pm 0.05$	$82.51 \pm 0.27$	<b><math>85.13 \pm 0.22</math></b>	$83.33 \pm 0.00$

## Contrastive Language-Image Pre-training (CLIP) [1]



Linear Probing: Fine-tuning the last linear layer only.



**Code:** <https://github.com/Haoxiang-Wang/ISR>

**Contact Information:**

- Haoxiang Wang: [hwang264@illinois.edu](mailto:hwang264@illinois.edu)
- Haozhe Si: [haozhes3@illinois.edu](mailto:haozhes3@illinois.edu)
- Bo Li: [lbo@illinois.edu](mailto:lbo@illinois.edu)
- Han Zhao: [hanzhao@illinois.edu](mailto:hanzhao@illinois.edu)