

# SEBRA : Debiasing through Self-Guided Bias Ranking



**ICLR**  
International Conference on  
Learning Representations



Adarsh Kappiyath<sup>1</sup>



Abhra Chaudhuri<sup>2</sup>



Ajay Kumar Jaiswal<sup>3</sup>



Ziquan Liu<sup>4</sup>



Yunpeng Li<sup>1</sup>



Xiatian Zhu<sup>1</sup>



Lu Yin<sup>1</sup>



<sup>1</sup>



<sup>2</sup>



**TEXAS**  
The University of Texas at Austin

<sup>3</sup>



<sup>4</sup>

# Preliminaries

Empirical Risk Minimization (ERM) with CE Loss exhibit tendency to learn different attributes asynchronously during training.

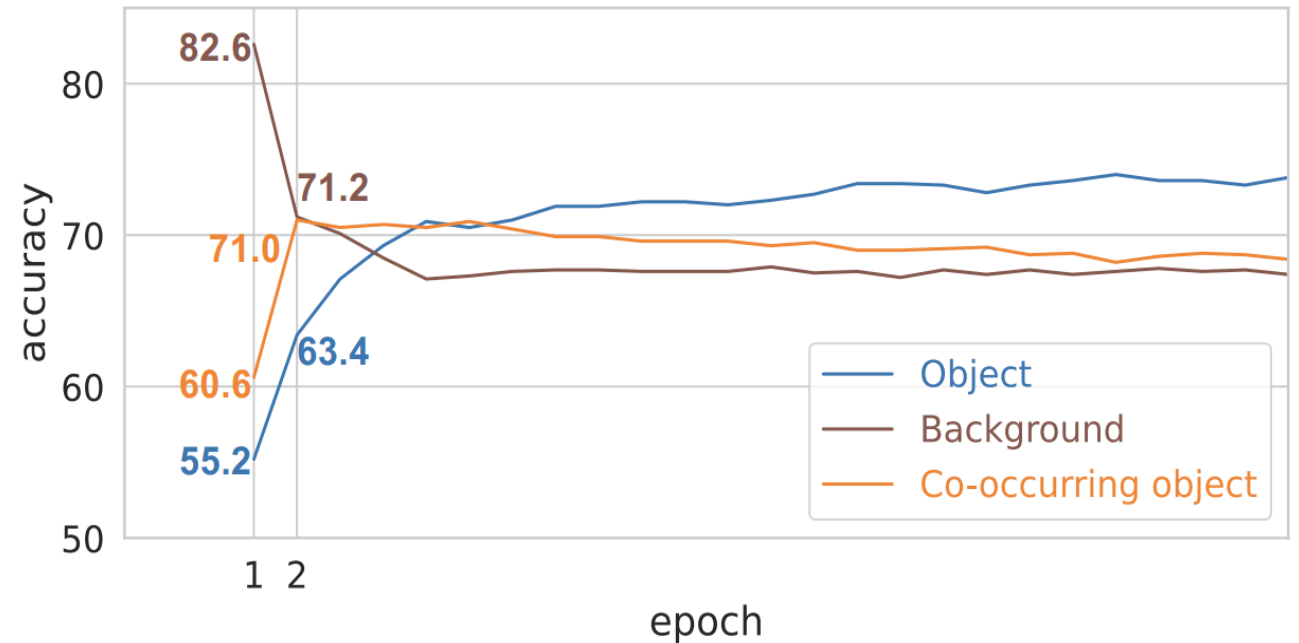


Fig 4 . Training Dynamics of Resnet50 with CE Loss on UrbanCars dataset.

# Preliminaries

Empirical Risk Minimization (ERM) with CE Loss exhibit tendency to learn different attributes asynchronously during training.

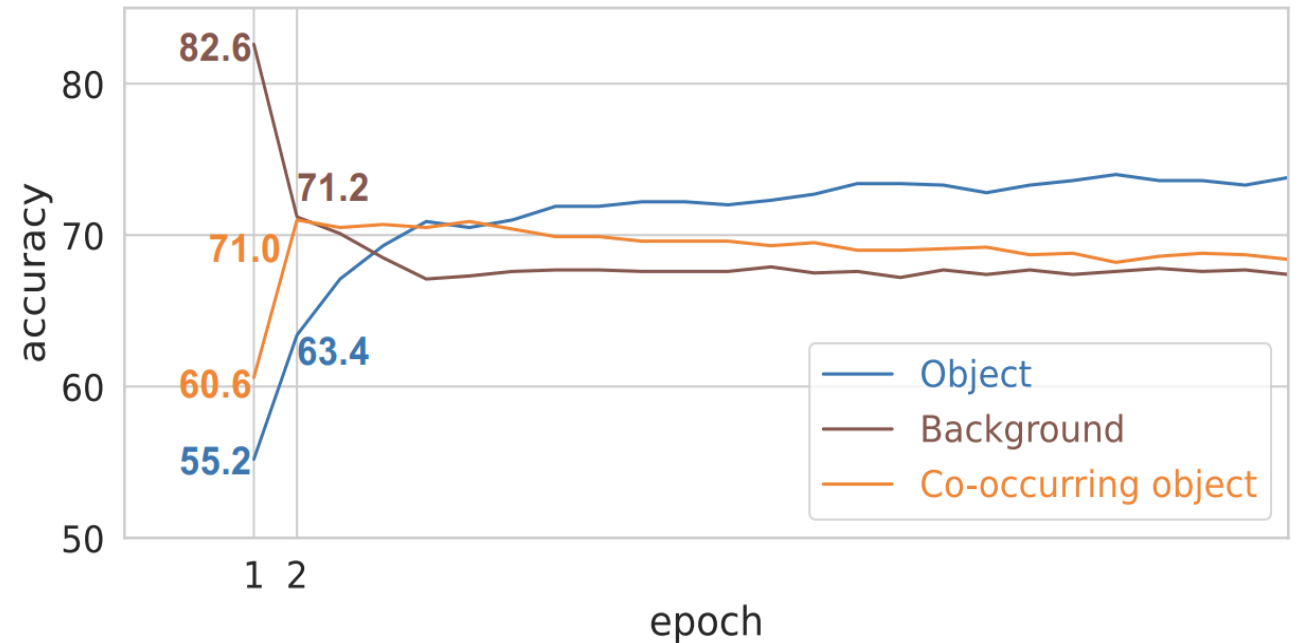


Fig 4 . Training Dynamics of Resnet50 with CE Loss on UrbanCars dataset.

*Idea: Modulate ERM dynamics to rank/order samples according to spuriousity. Mitigate biases based on ranking.*

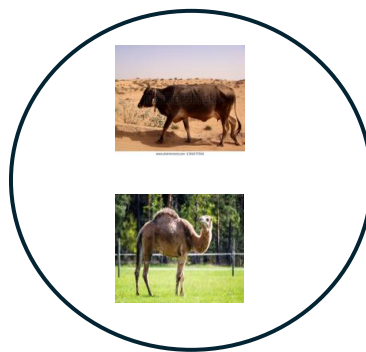
# Prior Works

Spurious Correlation  
Identification



Mitigating Impact of  
Spurious Correlations

- Identification of Spurious Correlations.
  - GCE Loss, Training with limited capacity models etc.



- Spuriousity of samples within a cluster
- Relative spuriousity across clusters



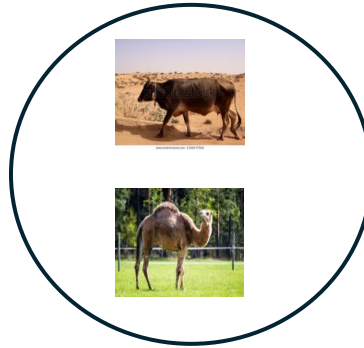
# Prior Works

Spurious Correlation  
Identification



Mitigating Impact of  
Spurious Correlations

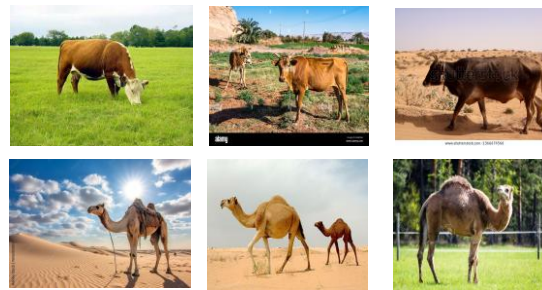
- Identification of Spurious Correlations.
  - GCE Loss, Training with limited capacity models etc.



- Spuriousity of samples within a cluster
- Relative spuriousity across clusters



- Class wise ranking of instances based on the *strength of spurious correlations(spuriousity)*.

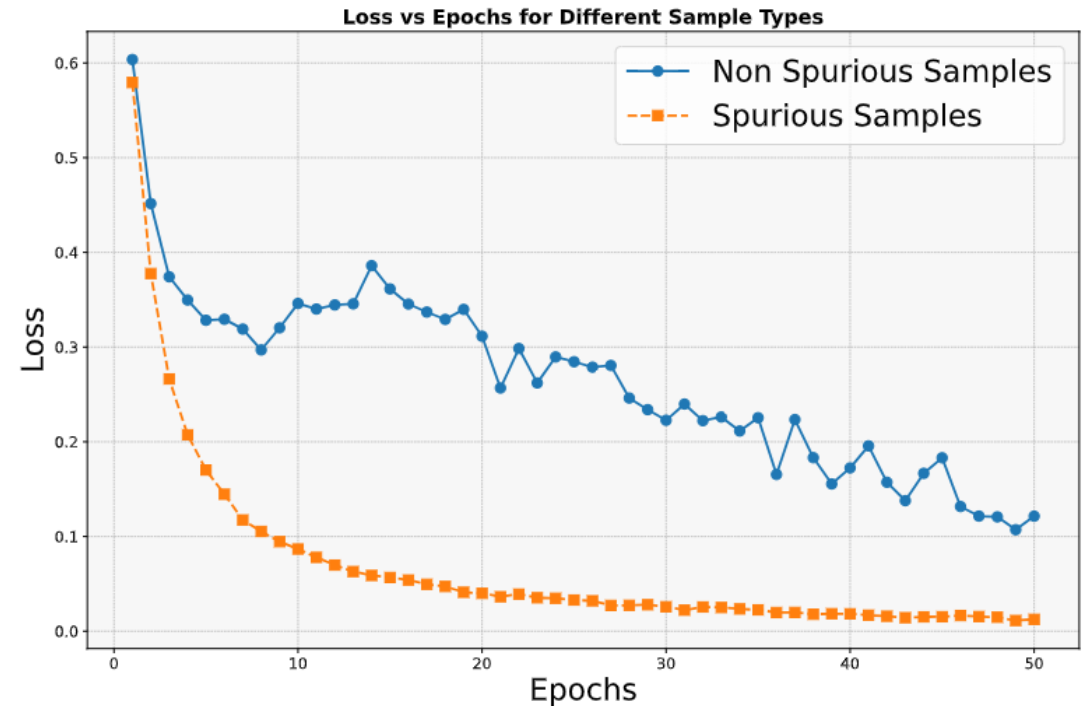
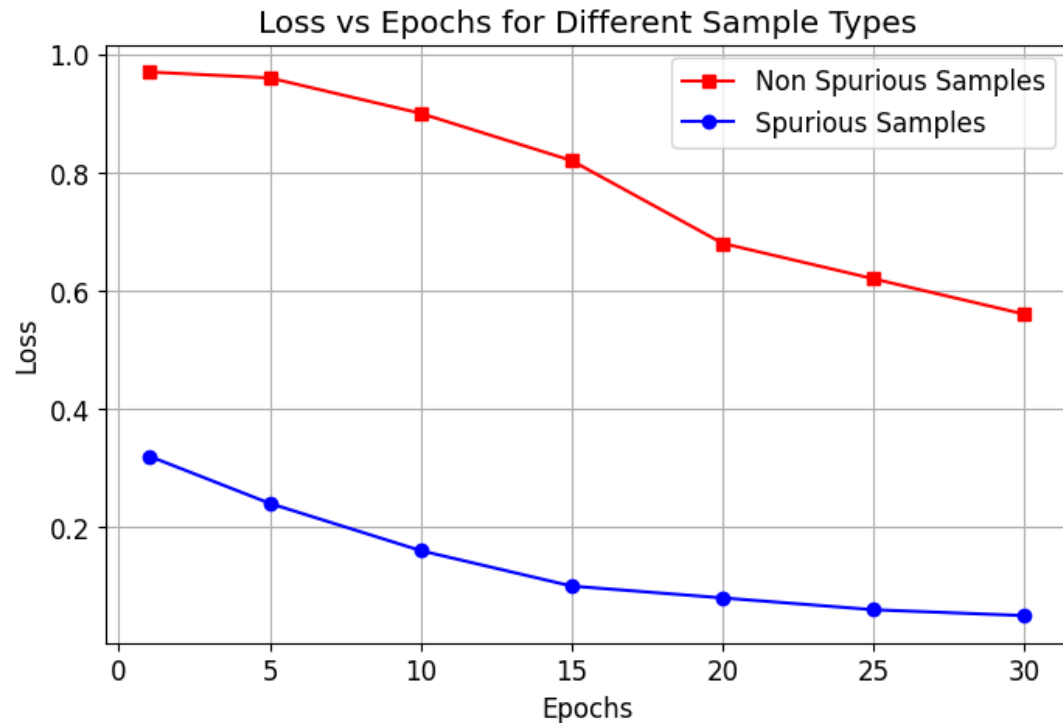


Human Supervision



# Assumption : Hardness Spuriousity Symmetry

The hardness of learning a sample, and its corresponding spuriousity measure, are symmetric to each other – the harder it is to learn a sample, the lower its spuriousity measure, and vice versa.



# Deviation of ERM in the Multi-Bias Setting

Global trends of ERM deviate due to:

- **Reliance on spurious features**
- **Non-uniform gradient updates.**

# Steering ERM in the Multi-Bias Setting

Correcting deviation requires explicit steering towards to maintain the Hardness-Spuriosity Symmetry.

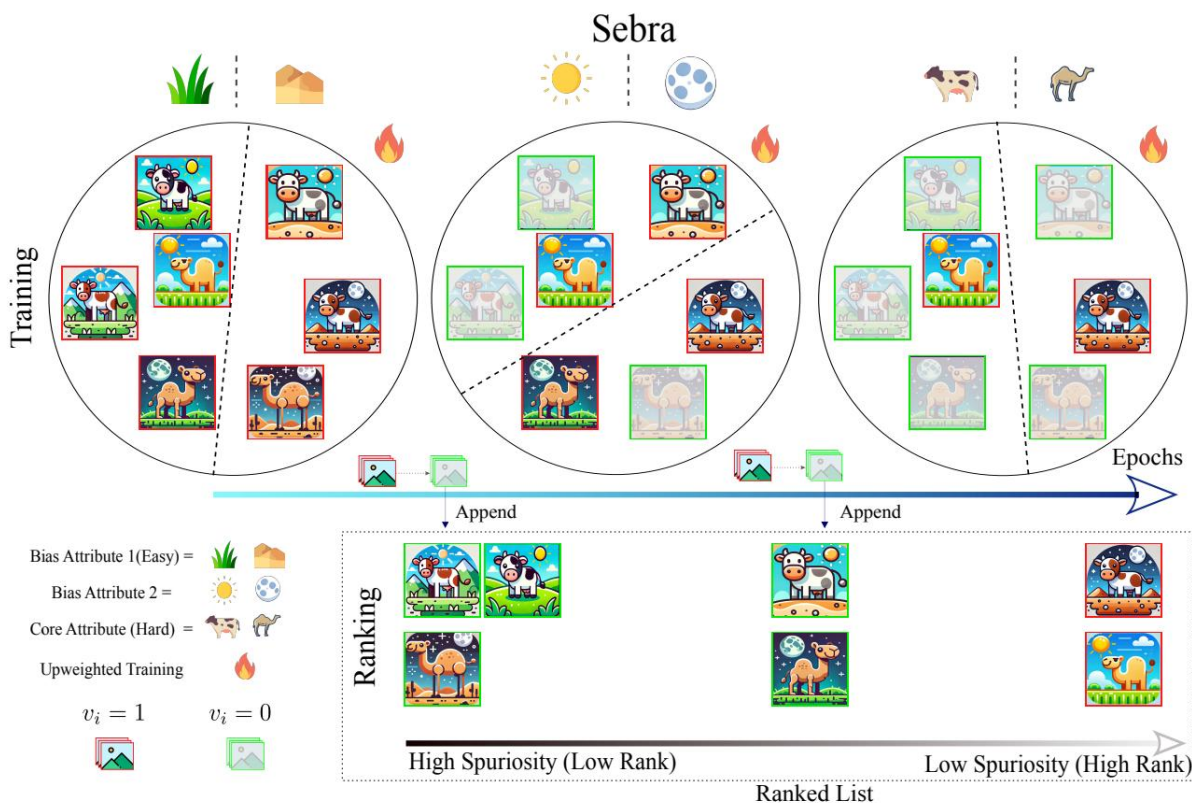


- One bias at a time



# Steering ERM in the Multi-Bias Setting

Correcting deviation requires explicit steering towards to maintain the Hardness-Spuriosity Symmetry.



- Spuriousity-Based Sequential Learning
- One bias at a time

# Sequential Learning Based on Levels of Spuriousity

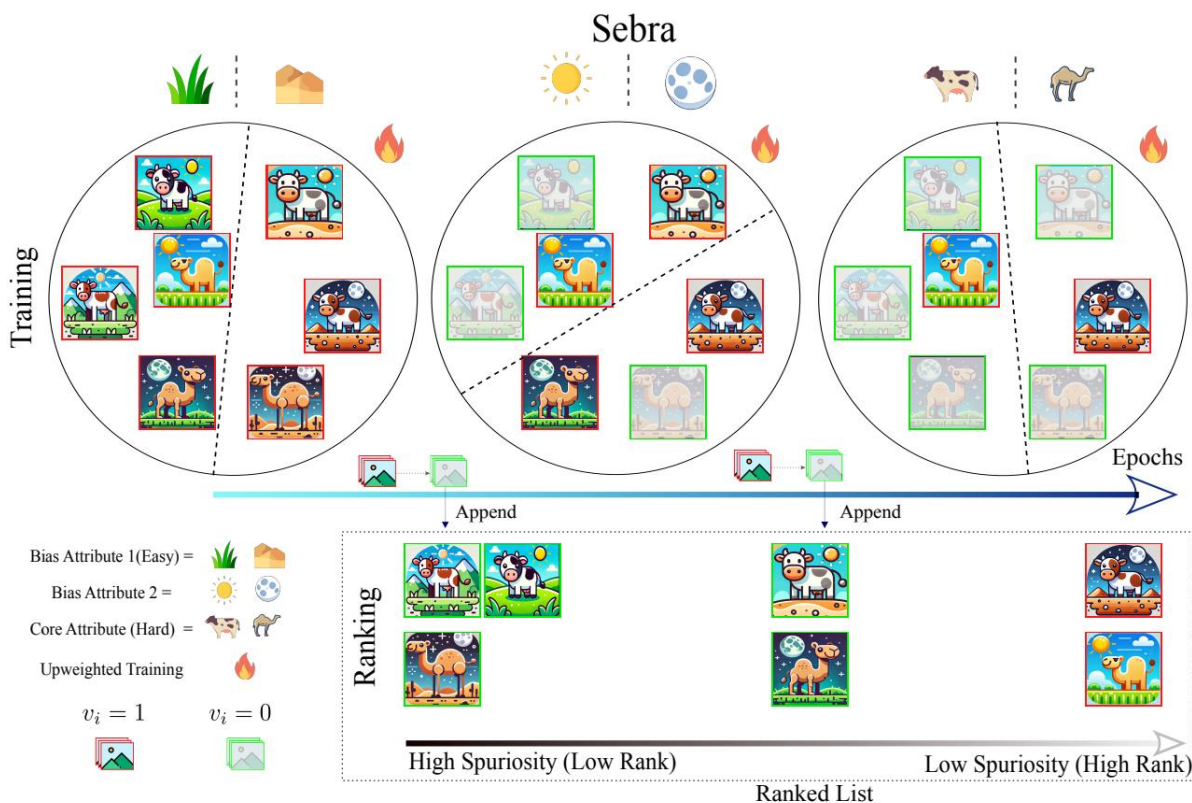
Selection variable ensures isolation among subgroups.



$$\max_v \sum_{i=1}^N \{ \underline{v_i^t} \mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i) - \lambda \underline{v_i^t} \}$$

# Sequential Learning Based on Levels of Spuriousity

Reweight CE-Loss by some measure of spuriousity. Vulnerable to shortcut  $u_i = 0$ , for all  $u_i$ s.

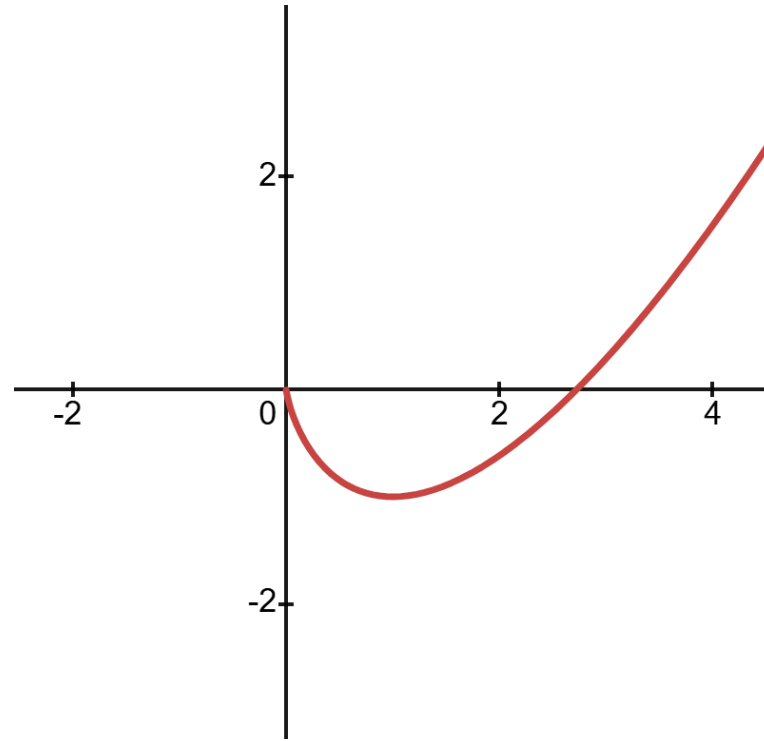


$$\min_{\theta} \min_u \sum_{i=1}^N \{ \underline{u_i} \mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i) \}$$

$$\max_v \sum_{i=1}^N \{ v_i^t \mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i) - \lambda v_i^t \}$$

# Hardness-Spuriosity Conservation Law

**Remedy for collapse:**  $u_i$ s must belong to a specific manifold satisfying a certain conservation law.

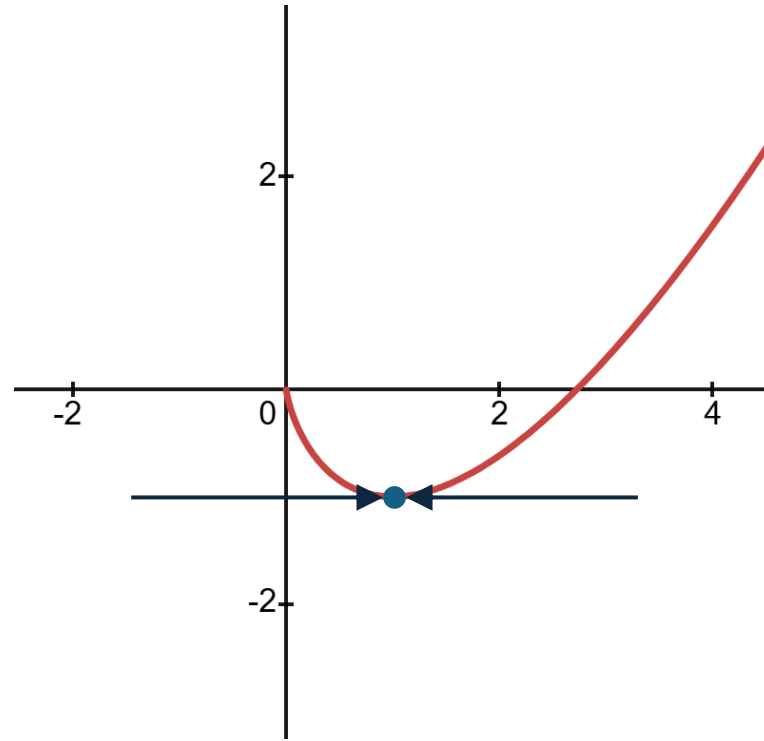


$$\min_{\theta} \min_u \sum_{i=1}^N \{ \underline{u_i} \mathcal{L}_{CE}(f_{\theta}(x_i), y_i) \}$$

$$u_i \mathcal{L}_{CE}(f_{\theta}(x_i), y_i) + \beta(u_i \ln u_i - u_i) = c$$

# Hardness-Spuriosity Conservation Law

**Remedy for collapse:**  $u_i$ s must belong to a specific manifold satisfying a certain conservation law. Ensure adherence to manifold through  $\beta$ -weighted regularization.



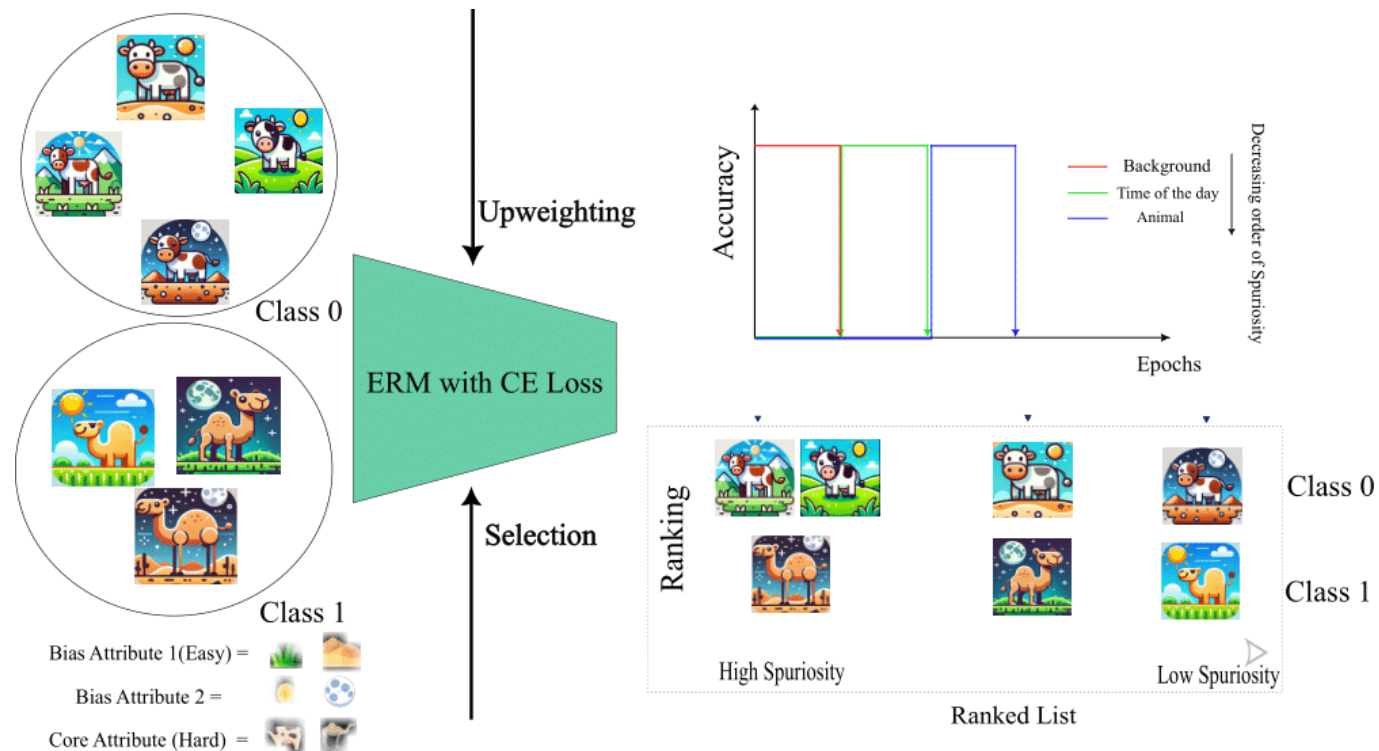
$$u_i \mathcal{L}_{CE}(f_\theta(x_i), y_i) + \beta(u_i \ln u_i - u_i) = c$$

$$\min_{\theta} \min_u \sum_{i=1}^N \{ \underline{u_i} \mathcal{L}_{CE}(f_\theta(x_i), y_i) \}$$

$$g(u_i) = (u_i \ln u_i - u_i)$$

$$\min_{\theta, u} \sum_{i=1}^N \{ \underline{u_i} \mathcal{L}_{CE}(f_\theta(x_i), y_i) + \beta \underline{g(u_i)} \}$$

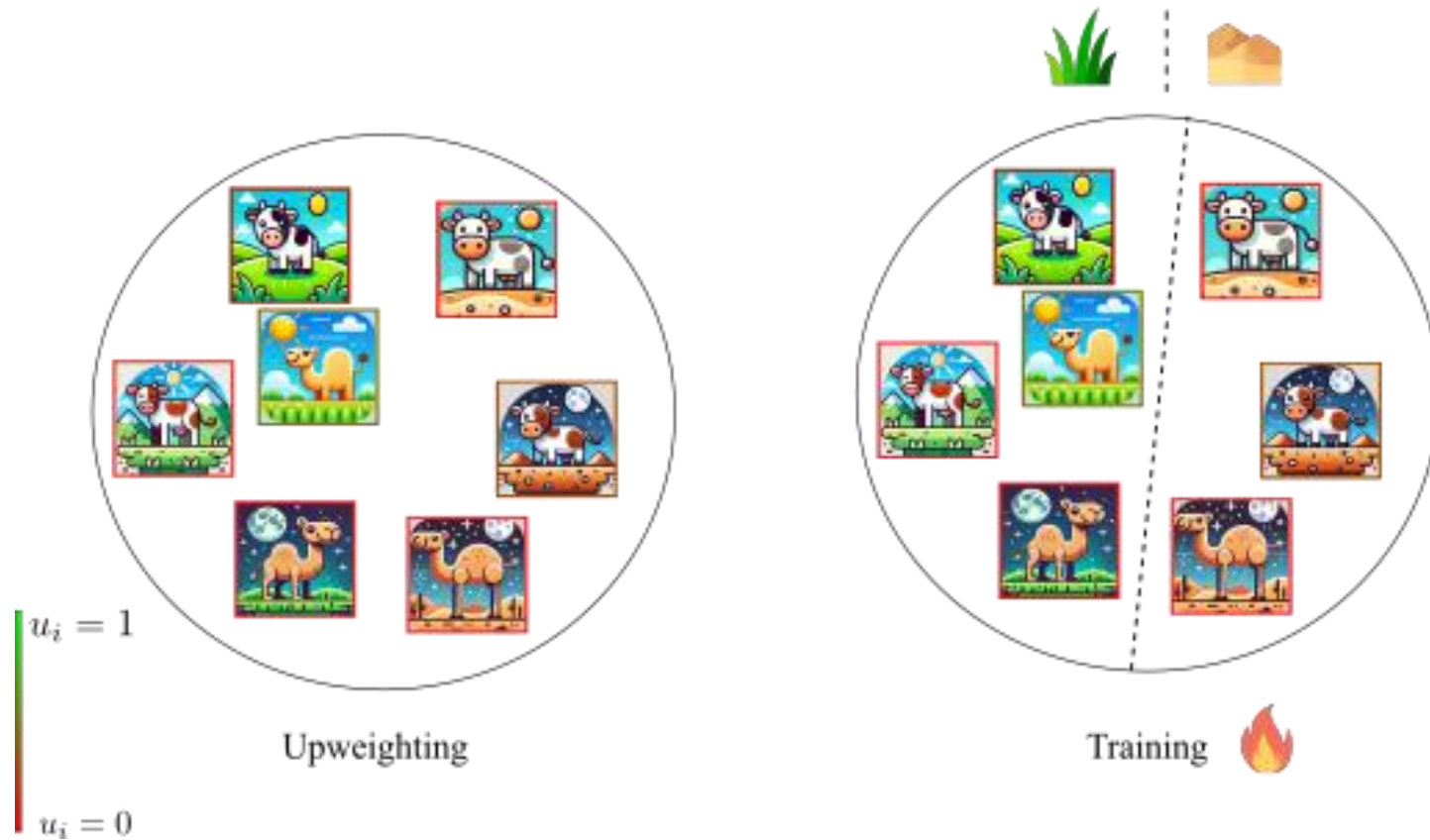
# Sebra : Self-Guided Bias Ranking



$$\mathcal{L}_{\text{ranking}}(\theta, u, v) = \sum_{i=1}^N v_i^{t-1} \{v_i^t u_i \mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i) - \lambda v_i^t - \beta u_i + \beta u_i \ln u_i\}$$

$$\min_{\theta, u} \max_v \mathcal{L}_{\text{ranking}}(\theta, u, v),$$

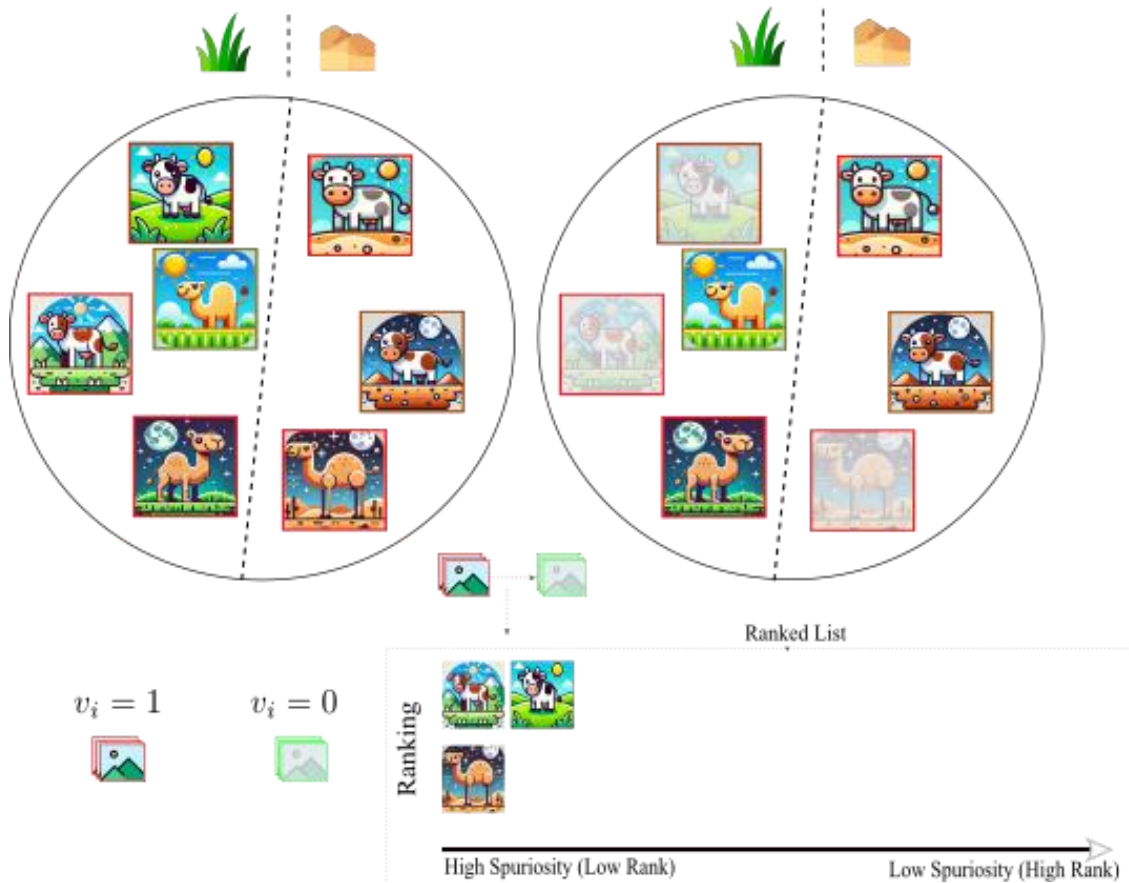
# Sebra : Upweighted Training



$$\min_{\theta, u} \sum_{i=1}^N \{ \underline{u_i} \mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i) + \beta g(\underline{u_i}) \}$$

$$u_i^* = p_y^{\frac{1}{\beta}}$$

# Sebra : Selection & Ranking

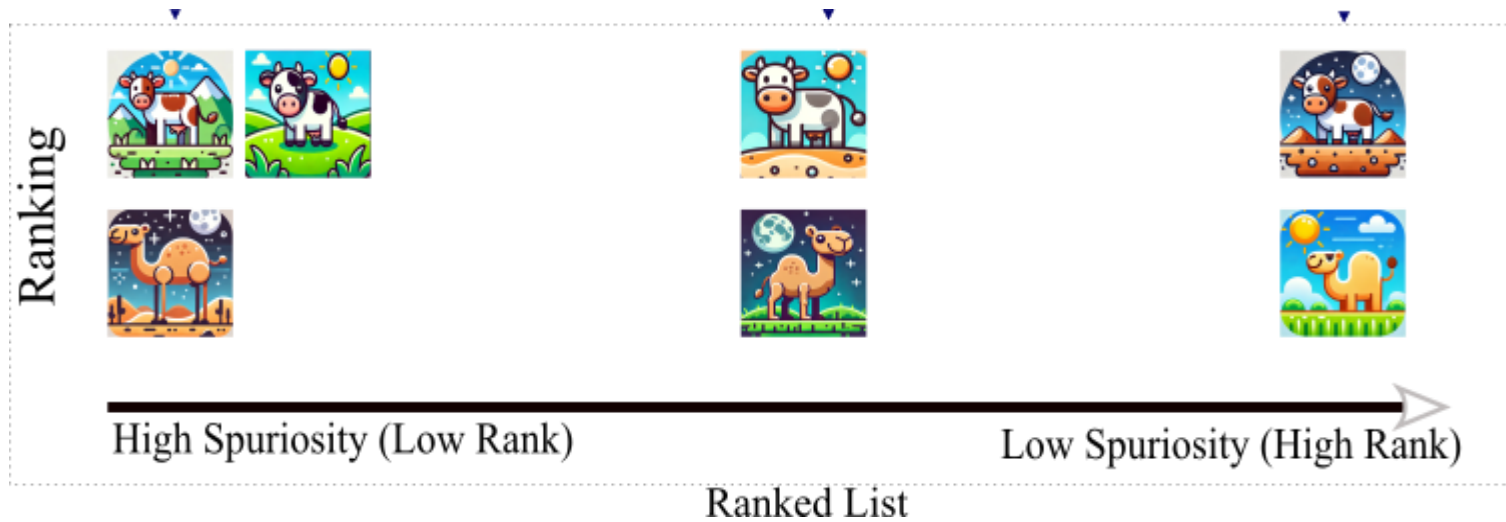


$$\max_v \sum_{i=1}^N \{ \underline{v_i^t} \mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i) - \lambda \underline{v_i^t} \}$$

$$v_i^{t*} = \begin{cases} 0, & \text{if } p_y > p_{\text{critical}}, \\ 1, & \text{otherwise.} \end{cases}$$



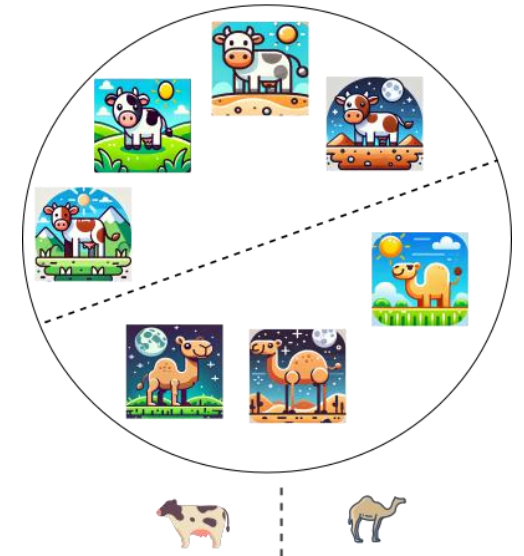
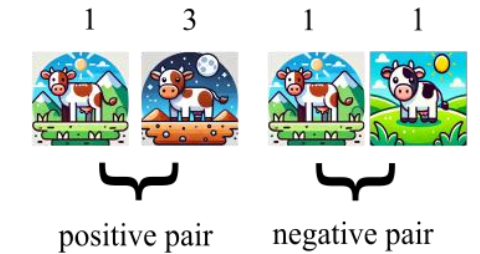
# Sebra : Debiasing



$$\mathcal{L}_{\text{con}}^{\text{sup}}(x; f_{\text{enc}}) = \mathbb{E} \left[ -\log \frac{\exp(z^{\top} z_m^{+} / \tau)}{\sum_{m=1}^M \exp(z^{\top} z_m^{+} / \tau) + \sum_{n=1}^N \exp(z^{\top} z_n^{-} / \tau)} \right],$$

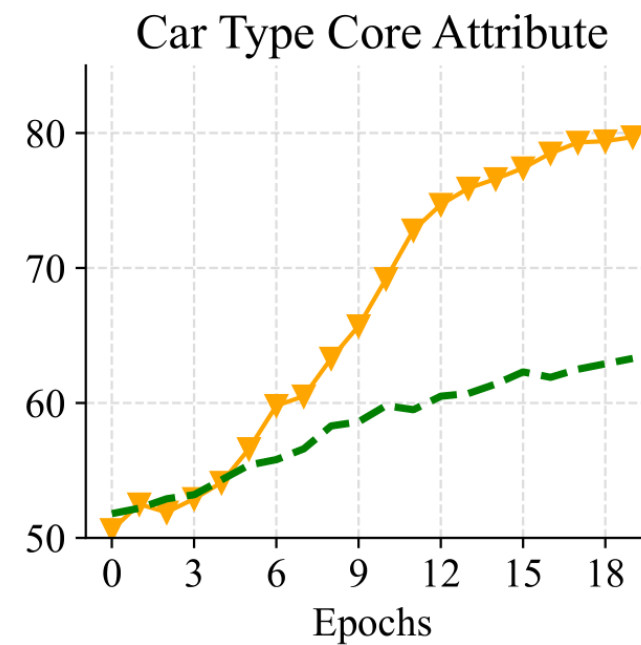
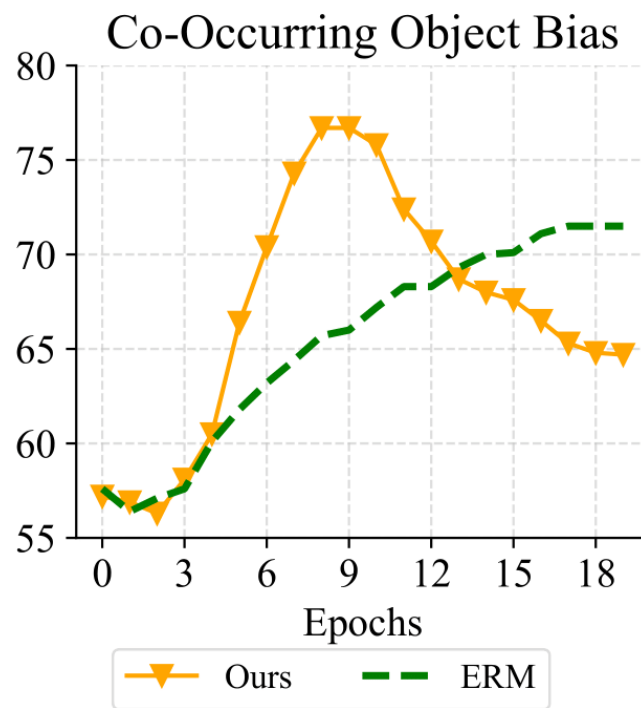
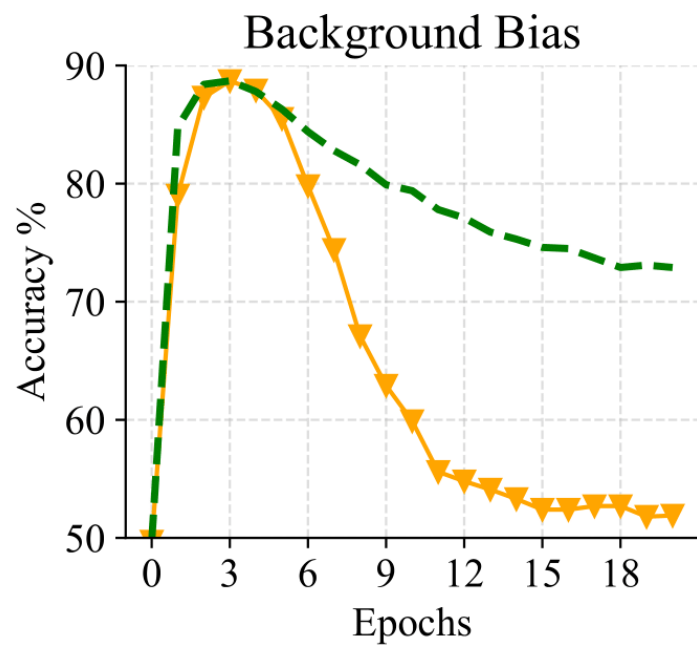
where  $\tau$  is the temperature coefficient,  $z_m^{+}$ ,  $z_n^{-}$  and  $z^{\top}$  are the embeddings of positive, negative, and reference samples respectively.

## Contrastive DeBiasing



# Modulation of ERM Dynamics

Sebra successfully mitigates the Whac-a-Mole Dilemma.



# Sebra : Ranking Results

Diving

Top Ranked



Bottom Ranked



Pole Vaulting

Top Ranked



Bottom Ranked



Table 1: Quantitative comparison of Sebra with various baselines. The results are shown in terms of Kendall's  $\tau$  for Urban Cars and CelebA, and Performance Disparity (PD) for BAR.

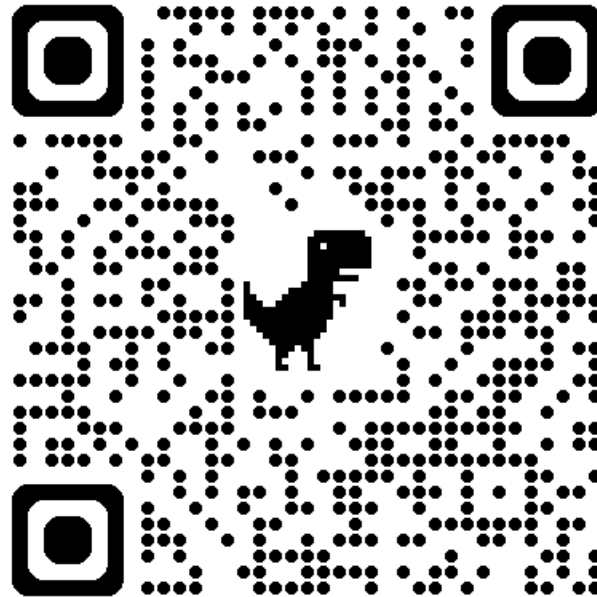
Method	Urban Cars	CelebA	BAR
Metric	Kendall's $\tau$ ( $\uparrow$ )	Kendall's $\tau$ ( $\uparrow$ )	PD ( $\uparrow$ )
Random Ordering	0.02	-0.01	0.25
ERM-based Ranking	0.12	0.14	4.55
Spuriousity Ranking	0.40	0.38	28.88
Sebra (Ours)	<b>0.85</b>	<b>0.69</b>	<b>32.32</b>

# Sebra : Debiasing Results

Methods	Sup.	UrbanCars			CelebA			BAR
		I.D. Acc. (↑)	WG Acc. (↑)	Avg GAP (↑)	I.D. Acc. (↑)	WG Acc. (↑)	Avg GAP (↑)	Test Acc. (↑)
Group DRO	✓	91.60 (1.23)	75.70 (1.79)	-10.30 (1.35)	90.08 (0.70)	37.9 (1.6)	-5.79 (1.63)	-
ERM	✗	97.60 (0.86)	33.20 (0.86)	-31.90 (3.92)	96.43 (0.13)	36.0 (1.7)	-22.83 (0.84)	68.00 (0.43)
LfF	✗	97.20 (2.40)	35.60 (2.40)	-31.06 (3.56)	95.12 (0.35)	35.5 (2.0)	-22.57 (1.26)	68.30 (0.97)
JTT	✗	95.80 (1.45)	33.30 (6.90)	-20.50 (2.61)	91.86 (1.48)	38.7 (2.4)	-26.81 (2.53)	68.14 (0.28)
Debian	✗	98.00 (0.89)	30.10 (0.89)	-31.40 (1.44)	96.28 (0.37)	41.1 (4.3)	-22.56 (0.54)	69.88 (2.92)
DFR	✗	89.70 (1.21)	-	-20.93 (2.61)	60.12 (1.28)	-	-19.16 (3.27)	69.22 (1.25)
Sebra (Ours)	✗	92.54 (2.10)	<b>73.8 (3.28)</b>	<b>-10.57 (1.72)</b>	88.61 (3.36)	<b>65.3 (4.1)</b>	<b>-9.82 (3.06)</b>	<b>75.36 (2.23)</b>

Method	Sup.	ImageNet-1K					MultiNLI
		I.D. Acc. (↑)	IN-W Gap (↑)	IN-9 Gap (↑)	IN-R Gap (↑)	Carton Gap (↑)	WG. Acc (↑)
LLE	✓	76.25	-6.18	-3.82	-54.89	+10	-
ERM	✗	76.13	-26.64	-5.53	-55.96	+40	66.8
LfF	✗	70.26	-17.57	-8.10	-56.54	+40	63.6
JTT	✗	75.64	-15.74	-6.75	-55.70	+32	69.1
Debian	✗	74.05	-20.00	-7.29	-56.70	+30	-
Sebra (Ours)	✗	74.89	<b>-14.77</b>	<b>-3.15</b>	<b>-54.81</b>	<b>+25</b>	<b>72.3</b>

# Thank You



[https://kadarsh22.github.io/sebra\\_iclr25/](https://kadarsh22.github.io/sebra_iclr25/)