

# SEBRA: DEBIASING THROUGH SELF-GUIDED BIAS RANKING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Ranking samples by fine-grained estimates of spuriousity (the degree to which spurious cues are present) has recently been shown to significantly benefit bias mitigation, over the traditional binary biased-*vs*-unbiased partitioning of train sets. However, this spuriousity ranking comes with the requirement of human supervision. In this paper, we propose a debiasing framework based on our novel Self-Guided Bias Ranking (*Sebra*), that mitigates biases via an automatic ranking of data points by spuriousity within their respective classes. Sebra leverages a key local symmetry in Empirical Risk Minimization (ERM) training – the ease of learning a sample via ERM inversely correlates with its spuriousity; the fewer spurious correlations a sample exhibits, the harder it is to learn, and vice versa. However, globally across iterations, ERM tends to deviate from this symmetry. Sebra dynamically steers ERM to correct this deviation, facilitating the sequential learning of attributes in increasing order of difficulty, *i.e.*, decreasing order of spuriousity. As a result, the sequence in which Sebra learns samples naturally provides spuriousity rankings. We use the resulting fine-grained bias characterization in a contrastive learning framework to mitigate biases from multiple sources. Extensive experiments show that Sebra consistently outperforms previous state-of-the-art unsupervised debiasing techniques across multiple standard benchmarks, including UrbanCars, BAR, and CelebA.

## 1 INTRODUCTION

Distribution shifts driven by spurious correlations (*aka* biases or shortcuts) are arguably one of the most studied forms of subpopulation shift (Koh et al., 2021; Yang et al., 2023). Models trained on data that have certain *easy-to-learn* attributes, spuriously correlated with labels, can overly rely on such spurious attributes, resulting in suboptimal performance during deployment (Geirhos et al., 2019). Both supervised (Sagawa et al., 2020; Idrissi et al., 2022) and unsupervised (Nam et al., 2020; Liu et al., 2021; Li et al., 2022; Park et al., 2023) methodologies for making neural networks robust to spurious correlations, a task also known as *debiasing*, have been developed. To get around the expensive human labor involved in acquiring bias labels for training supervised debiasing algorithms, unsupervised methods typically take a two-stage approach: an initial stage for bias identification and a second stage for bias mitigation. Unsupervised bias identification often relies on certain characteristics of spurious attributes, such as their relative ease of learning compared to target attributes (Nam et al., 2020), formation of clusters in feature space (Sohoni et al., 2020), adherence to a low-rank property (Huh et al., 2023), etc. This is followed by the mitigation step via resampling (Idrissi et al., 2022), contrastive learning (Zhang et al., 2022), and pruning (Park et al., 2023), etc.

Existing bias identification methods typically categorize data points into two (Nam et al., 2020; Liu et al., 2021) or more discrete groups (Sohoni et al., 2020; Yang et al., 2024). However, they do not offer insights into how the strength of spurious correlations varies across the identified groups, nor do they account for the variation in the strength of spurious correlations across instances within each group. Recent works, such as Singla & Feizi (2022); Moayeri et al. (2023), address these limitations by ranking data points based on spuriousity—the degree to which common spurious cues are present. However, these [\[rebuttal: spuriousity / bias ranking\]](#) methods rely on human supervision or auxiliary biased models to identify biased features. Furthermore, they heavily rely on the interpretability of neural features extracted from adversarially trained encoders and the effectiveness of the interpretability techniques employed, which limits their applicability.

To address these limitations, we present Self-Guided Bias Ranking (*Sebra*), a spuriousity ranking algorithm without the need for human intervention. *Sebra* is based on the observation of a local symmetry in ERM (Empirical Risk Minimization) training – in a given iteration, the hardness of learning a sample is inversely correlated with the amount of spurious features it contains. In other words, the lower the amount of spurious features, the harder a sample is to learn, and vice versa. [rebuttal: We call this, the Hardness-Spuriosity Symmetry (Assumption 1), which consequently gives rise to a corresponding conservation law (Theorem 1) relating the hardness of learning a sample to a measure of its spuriousity. This implies that the spuriousity ranking can be derived by looking at the trajectory (through the sample space) of a model that learns attributes sequentially in increasing order of hardness.]

[rebuttal: However, when training a neural network on samples with varying levels of spuriousity, globally across iterations, ERM tends to deviate from this trajectory due to (a) reliance on spurious features, since higher spuriousity samples are known to inhibit the learning of those with relatively lower levels of spuriousity Qiu et al. (2024) and (b) non-uniform gradient updates received for samples of different levels of spuriousity due to different values of the task loss (influenced by their levels of spuriousity), leading to non-determinism in the order in which samples are learned. *Sebra* corrects this deviation by steering the optimization pathway of a neural network by dynamically modulating ERM through a pair of controller variables to follow the conservation law corresponding to the Hardness-Spuriosity Symmetry, while minimizing the interference caused by samples of one spuriousity level on the learning of the other, thereby enabling the network to learn attributes in increasing order of difficulty.] Consequently, a readout of the order in which samples are learned along this pathway serves as our predicted spuriousity ranking, requiring no human supervision. By leveraging the fine-grained spuriousity rankings obtained through *Sebra* and incorporating them into a simple contrastive loss, we outperform previous state-of-the-art unsupervised and supervised debiasing techniques across multiple standard benchmarks, including UrbanCars, BAR, and CelebA.

To summarize, we: ❶ introduce a novel self-guided bias ranking framework, *Sebra*, to dynamically rank the data points of each class on the decreasing order of the strength of spurious signals, without any human supervision; ❷ utilize these derived rankings to enable debiased learning using a simple contrastive learning framework; ❸ empirically demonstrate the effectiveness of our proposed approach across multiple datasets with spurious correlations, including UrbanCars, CelebA, and BAR. Our method achieves an average improvement of 10% in both UrbanCars and CelebA, and 6% in BAR under subpopulation shift, outperforming state-of-the-art unsupervised debiasing approaches.

## 2 RELATED WORKS

**Bias Identification:** A plethora of methods assume knowledge of bias either in the form of bias labels Lee et al. (2021); Idrissi et al. (2022) or type of bias Geirhos et al. (2019); Chang et al. (2021). Even though these methods produce superior debiasing results, obtaining bias annotations for all biases or identifying the type of bias requires significant human efforts. This led to the development of various inductive biases suitable for bias identification. One of the most commonly used inductive biases for bias identification is the property of bias being easier to learn. In Nam et al. (2020), bias is identified by obtaining a bias-only model through upweighting data points that are easy to learn. Another popular bias identification strategy relies on training a model with a limited network capacity using empirical risk minimization Liu et al. (2021), the hypothesis being that a model with a small capacity would face difficulties in learning complex features and thus prefer to learn easy spurious features. Such simple bias identification schemes has shown to be very useful for datasets with single bias attributes but encounter *Whac-A-Mole dilemma* Li et al. (2023) when faced with datasets with multiple spurious correlations. In Sohoni et al. (2020); Yang et al. (2024), clusters based on biased attributes in the feature space are utilized for bias identification but provide no means to characterize the nature of clusters discovered or how the strength of spurious attributes varies across these discovered clusters. Recently, Singla & Feizi (2022) proposed a method to identify spurious and core attributes by analyzing neural features of adversarially trained encoders using interpretability techniques like GradCAM and feature attacks. Building on these insights, Moayeri et al. (2023) rank instances within a class based on the presence of these identified attributes, sorting data in decreasing order of spuriousity. Although these methods offer a detailed characterization of spurious attributes in the dataset, their dependence on human supervision and the quality of interpretability techniques used can restrict their applicability. In contrast, our proposed ranking

framework orders data in decreasing order of easiness to learn as perceived by an ERM model, without relying on human supervision or fragile interpretability techniques.

**Bias Mitigation:** Some of the simple bias mitigation strategies involve up-weighting bias conflicting points and down-weighting bias aligned points, thereby promoting the model to learn target features from the data Liu et al. (2021); Idrissi et al. (2022); Lee et al. (2021). Other approaches include obtaining a debiased model by training a model to learn different mechanisms to that of a bias-only model Nam et al. (2020), pruning Park et al. (2023) or forgetting the bias information from a biased model Tiwari & Shenoy (2023). In the presence of group labels either inferred or via supervision, a debiased model is obtained by minimizing worst group risk Sagawa et al. (2020). Although simple upweighting methods have shown to be very effective in debiasing they lead to the underutilization of diversity of the training data resulting in suboptimal performance. With the more fine-grained bias identification scheme, we utilise the available data more efficiently using contrastive loss to facilitate debiasing. Contrastive learning effectively debiases data Zhang et al. (2022); Jung et al. (2023), but our ranking scheme refines pair selection, boosting debiasing performance and scalability to diverse and large-scale datasets.

### 3 METHODOLOGY

In this section, we introduce a novel spuriousity-ranking framework, Sebra, designed to rank or order data points in decreasing order of spuriousity. At its core, the framework integrates self-guided weighting mechanisms into the standard Empirical Risk Minimization (ERM) using cross-entropy loss, creating an objective that prioritizes data points by their spuriousness. These self-guided weighting mechanisms guide ERM consistently along a pathway wherein attributes are learned sequentially in the increasing order of hardness. As a result, the order in which instances transition from unlearned to learned naturally reflects the spuriousity of the data point. We demonstrate the effectiveness of this ranked dataset for debiasing within a contrastive learning framework. Our approach is formalized in Section 3.1, with a diagrammatic illustration in Fig. 1.

#### 3.1 SEBRA: SELF-GUIDED BIAS RANKING

**Intuition behind Sebra:** [rebuttal: Following the example in Fig. 1, consider the problem of classifying “cows” and “camels”, where in the train set, cows are spuriously correlated with “green” backgrounds (such as grasslands) in “daylight”, and camels are spuriously correlated with the “desert” background at “nighttime”. Now, a model trained with ERM tends to classify the training datapoints first based on the background, *i.e.*, cows on grasslands *v.s.* camels on deserts, which implies that it is the easiest attribute to learn Nam et al. (2020). However, when samples exhibiting the background spurious correlation are dropped out from the training set, ERM learns to classify based on the lighting conditions, *i.e.*, cows in daylight *v.s.* camels at nighttime. Finally, it is only when these are also dropped from the training set does the model finally capture the core attributes of cows and camels. Thus, when controlled with an appropriate steering mechanism (dropping of training samples corresponding to the already-learned spurious attribute), the sequence in which ERM learns data points follows a *high-spuriousity to low-spuriousity pathway*, naturally providing a spuriousity ranking. This fine-grained ordering can then be exploited through contrastive learning for debiasing.]

**Notations:** Given a train set  $X = \{(x_i, y_i)\}_{i=1}^N$  with  $N$  data points across  $C$  classes, we aim to rank them in decreasing order of spuriousity, *i.e.*, if  $x_i$  exhibits spurious cues than  $x_j$ , then  $\rho(x_i) < \rho(x_j)$ , where  $\rho(x)$  is an integer in  $[0, N]$  indicating the spuriousity rank of  $x$ . We use a neural network  $f_\theta$  with parameters  $\theta$  to drive the ranking process. All proofs and derivations are provided in Appendix A.

**Definition 1.** For a sample  $x \in X$ , let  $F_x$  be the set of all features types / attributes in  $x$ . An attribute space  $\mathcal{A}$  is the exhaustive collection of all feature types across all  $x \in X$ , *i.e.*,

$$\mathcal{A} = \bigcup_{x \in X} F_x$$

**Definition 2 (Attribute Types and Spuriousity Ranking).** A *causal attribute*  $a_c \in \mathcal{A}$  is one that is responsible for determining the label  $y$  of a datapoint  $x, \forall x \in X$ . A *spurious attribute*  $a_s \in \mathcal{A}$  is a non-causal attribute that does not determine the label  $y$  of any sample  $x \in X$ , but co-occurs frequently with  $a_c$  in  $X$ . We call the subspace of  $\mathcal{A}$  covering all spurious features,  $\mathcal{A}_s$ , the *spuriousity basis*. The *spuriousity measure*  $\mu(x)$  on  $X$  is the fraction of spurious attributes in  $\mathcal{A}_s$  spanned by the

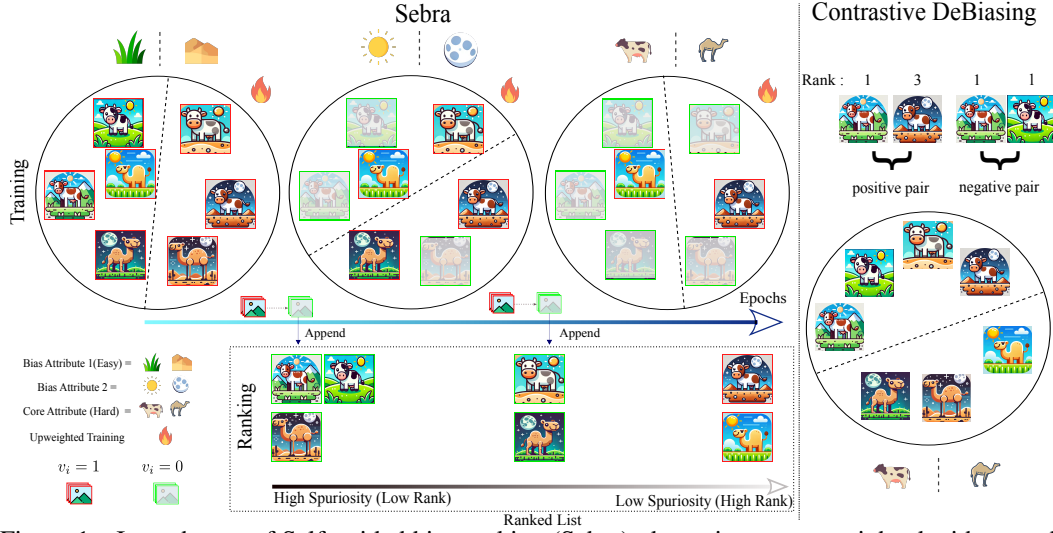


Figure 1: In each step of Self-guided bias ranking (Sebra), datapoints are upweighted with  $u_i$  and then trained via ERM. Following this, we estimate  $v_i$  for each sample to select them for subsequent training. Samples for which  $v_i$  transitions from 1 to 0 are ranked at each step and eliminated from subsequent training. Any unranked samples are appended to the ranked list at the end of the training phase. In the mitigation phase, negative pairs are formed using samples with the same rank, while positive pairs are obtained using samples with a higher rank than the reference samples.

feature-type set  $F_x$  of a sample  $x \in X$ . A *spuriousity ranking*  $\rho(x)$  is an ordering on  $X$  such that:

$$\rho(x_i) < \rho(x_j) \quad \forall i, j \leq N \mid \mu(x_i) > \mu(x_j)$$

In other words, samples with high levels of spuriousity appear earlier in the ranking via  $\rho$  than samples with lower levels of spuriousity.

**Assumption 1** (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.

**Implementation of  $\rho(x)$ :** We leverage the Hardness-Spuriousity Symmetry to design a form of self-guided bias identification, steering ERM consistently along a *high-spuriousity to low-spuriousity pathway*. This results in the rank of a data point  $x_i$  being the epoch in which its cross-entropy loss (or a monotonically increasing function of it) drops below a certain threshold, or in other words, its predicted probability  $p_y$  (or a monotonically increasing function of it) of the correct class  $y$  exceeds a certain threshold, determined by hyperparameters, as discussed below.

**Fine-Grained Rank Resolution:** Note that this specific criterion of ranking maps the  $N$  datapoint to  $M$  buckets, where  $M \leq N$ . In other words, multiple datapoints can get mapped to the same rank bucket, if they transition below the loss / probability threshold together in the same iteration. However, in our implementation, we also provide information about the spuriousity measure  $\mu(x)$  for every data point  $x_i$  through a weighting factor called  $u_i \propto \mu(x)$ . Since  $\mu(x)$  is a continuous-valued function, sorting in the decreasing order of  $\mu(x)$  provides a straightforward mechanism for collision resolution and obtaining a fine-grained ranking among data points that inhabit the same coarse-grained rank bucket.

### 3.1.1 FORMULATION

Our ranking algorithm involves the following three key phases in each epoch:

- ① **Selection:** [rebuttal: We design the selection mechanism to shift the model’s focus to a new subgroup once a particular subgroup has been learned, for it to capture attribute types in order of increasing difficulty across iterations. The training set is partitioned into samples that have been learned, *i.e.*, easier samples with high spuriousity, and those that have not yet been learned, *i.e.*, difficult samples with low spuriousity. The latter are carried forward for further updates to  $\theta$  via



ERM. This segregation-based selection serves a dual purpose: it mitigates the influence of highly spurious features on the learning of the less spurious ones, and it promotes the learning of attributes in increasing order of difficulty.] To implement this, we introduce a binary selection variable  $v_i$  for each point  $x_i$ , which identifies a minimal subset of data points that maximizes the cross-entropy loss:

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

$v_i^t \mathcal{L}_{CE}(f_{\theta}(x_i), y_i)$  is responsible for selecting points that have not yet been learned (*i.e.*, those with a high  $\mathcal{L}_{CE}$ ), while the  $-\lambda v_i$  prevents the trivial solution where all  $v_i^t$ s are set to 1, minimizing the number of points that are selected in a single epoch.

Furthermore, to mitigate the influence of previously learned highly spurious attributes on subsequent learning, we condition the optimization on the state of  $v_i$  in the previous iteration, *i.e.*, on  $v_i^{t-1}$  as follows:

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

and additionally restrict the domain of  $v_i^t$  to  $\{0, v_i^{t-1}\}$  (instead of the general binary  $\{0, 1\}$ ), where  $v_i^0 = 1, \forall i \in [1, N]$ . This dynamic domain constraint follows from the order on  $X$  induced by the measure  $\mu$ . It effectively implements the inductive bias that points with higher bias would always be learned before points with fewer spurious features, leading to the result that once something has been learned and ranked (with their corresponding  $v_i^t$  set to 0), they need not be considered anymore. Note, however, that before  $v_i^{(t-1)}$  becomes 0, (*i.e.*, while it is still 1), solving for the optimal  $v_i^t$  is still effectively a general binary optimization problem on  $\{0, 1\}$ .

② **Upweighting:** Next, to counteract the non-uniform gradient updates inherent in ERM, and to facilitate the ranking of points with high spuriousity before those with low spuriousity, we utilize the inductive bias that ERM has a lower local risk, in any given iteration, for high spuriousity samples relative to their lower spuriousity counterparts (Assumption 1). We do so by introducing a weighting variable  $u_i$  for each point  $x_i$  proportional to the value of the spuriousity measure for that point,  $\mu(x_i)$  as follows:

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

This essentially results in the selection of those points with the lowest  $\mathcal{L}_{CE}$  and having them ranked before any of the other points with higher values of  $\mathcal{L}_{CE}$  (more difficult-to-learn points, and hence, with fewer spurious features). In principle, following from Assumption 1,  $u$  could be any monotonically decreasing function of  $\mathcal{L}_{CE}$ .

However, a shortcut solution to minimizing  $u$  is to set all  $u_i = 0$ . We prevent this shortcut by incorporating the inductive bias that  $u_i \propto \mu(x_i)$  into the optimization objective. For our specific case, we use  $u_i = e^{-t(\mathcal{L}_{CE}(f_{\theta}(x_i), y_i))}$ , which has the effect that samples with high learnability / spuriousity are upweighted by an exponential function of their spuriousity, where  $t(x)$  is a monotonically increasing function of  $x$ . We incorporate this constraint into the objective as follows:

$$\min_{\theta, u} \max_v \sum_{i=1}^N v_i^{t-1} \{v_i^t u_i \mathcal{L}_{CE}(f_{\theta}(x_i), y_i) - \lambda v_i^t + \beta g(u_i)\}, \quad (1)$$

where  $\beta$  is a hyperparameter determining the weight of this constraint, and  $g(u_i)$  is a convex function meant to impose the constraint, whose form we uncover next.

**Theorem 1** (Hardness-Spuriosity Conservation). *Iff the spuriousity measure  $u_i^* = e^{-t(\mathcal{L}_{CE}(f_{\theta}(x_i), y_i))}$ , where  $t(x)$  is a monotonically increasing function of  $x$ , the variable  $u_i$  in Eq. (1), across all values of  $\mathcal{L}_{CE}(f_{\theta}(x_i), y_i)$ , satisfies the following conservation law:*

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

such that  $u_i^*$  is the minimizer of the conserved function.

*Intuition:* Theorem 1 arises as a consequence of the Hardness-Spuriosity Symmetry (Assumption 1), which requires the measures of hardness ( $\mathcal{L}_{\text{CE}}$ ) and spuriousity ( $u_i$ ) to balance each other out. It states that, for the solution to Eq. (1) to have the form  $e^{-t(\mathcal{L}_{\text{CE}}(f_\theta(x_i), y_i))}$ , the quantity  $u_i \mathcal{L}_{\text{CE}}(f_\theta(x_i), y_i) + \beta(u_i \ln u_i - u_i)$  should be conserved, *i.e.*, a constant, for all valid choices of  $u_i$ . The implication is that the optimization on  $u$  should be restricted to the space of those values that follow the conservation law. It formalizes the constraint that we need to impose on  $u$  in order to avoid the shortcut of setting all  $u_i = 0$ .

In other words, the solution to  $u$  in Eq. (1) is the minimum in the space of all values that satisfy the conservation law. Based on this, we use  $g(u_i) = \beta(u_i \ln u_i - u_i)$  in Eq. (1) to enforce the conservation criterion, and obtain our final objective, which we optimize for all three sets of variables  $\theta$ ,  $u$ , and  $v$ :

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

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

③ **Ranking:** [rebuttal: Finally, samples with high spuriousity, *i.e.*, the ones that have been already learned and dropped out of the training set in the selection phase, are appended to the rank list.] Specifically, in every epoch  $t$ , we select those  $x_i$ s for which  $v_i = 0$  from the selection step, and append them to a rank list  $X_{\text{ranked}}$  (which is initially empty) as:

$$X_{\text{ranked}}^t = X_{\text{ranked}}^{t-1} || R,$$

where  $||$  is the concatenation operator between two lists,  $R$  is an ordered list of data points  $x$  such that  $x_i < x_j \implies u(x_i) \geq u(x_j)$ ;  $\forall x_i, x_j \in R$  and  $v^{t-1}(x) = 1, v^t(x) = 0$ ;  $\forall x \in R$ . [rebuttal: Below we discuss how Sebra progressively orders training samples based on spuriousity by optimizing the variables associated with the above three phases.]

### 3.2 OPTIMIZATION

[rebuttal: Based on our formulation in Section 3.1.1, Sebra is parameterized by a set of three variables,  $\theta$ ,  $u$ , and  $v$ , respectively corresponding to the Selection, Upweighting, and Ranking phases.] Since they are all independent, one can optimize  $\mathcal{L}_{\text{ranking}}$  wrt each of the variables by keeping the others fixed. In each iteration, we first solve for  $v_i^t$  to select the points that have not yet been sufficiently learnt, compute their corresponding  $u_i$ s, with which we upweight and minimize  $\mathcal{L}_{\text{CE}}$  wrt  $\theta$ , (which is non-zero for only those samples that have been selected by  $v$  in the beginning of the iteration), and finally, set aside and rank samples whose  $v_i$ s switched from 1 to 0 in this iteration to avoid interfering with subsequent rankings.

**Selection:** We start by maximizing  $\mathcal{L}_{\text{ranking}}(\theta, u, v)$  wrt  $v$ . Note, here, that solving for  $v_i^t$  is a discrete optimization problem, since  $v_i^t \in \{0, v_i^{t-1}\}$ . Let  $k = u_i \mathcal{L}_{\text{CE}}(f_\theta(x_i), y_i) - \lambda$ . This partitions the search space into two halves, *i.e.*,  $k \geq 0$  and  $k < 0$ , as follows:

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

The optimal  $v_i$  can be obtained in terms of the predicted probability of the correct class  $p_y$  as:

$$v_i^{t*} = \begin{cases} 0, & \text{if } p_y > e^{-\lambda/u_i}, \\ 1, & \text{otherwise.} \end{cases} \quad (2)$$

Once the  $v_i^t$  for a data point  $x_i$  has been set to 0, we consider it as learned, and by the design of the optimization objective, it does not influence the subsequent learning of the remaining points.

**Upweighting and Training:** We then solve for the minimization  $\mathcal{L}_{\text{ranking}}(\theta, u, v)$  in  $u$ :

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

Since  $\mathcal{L}(u \mid \theta, v)$  is a convex function, it can be minimized by equating its derivative wrt  $u$  to 0 and solving the resulting equation (when  $v_i^t = 1$ ), which gives us the minimizer of  $\mathcal{L}(u \mid \theta, v)$  as:

$$u_i^* = p_y^{\frac{1}{\beta}} \quad (3)$$

We then optimize the parameters of the neural network  $\theta$  as follows via regular mini-batch stochastic gradient descent:

$$\min_{\theta} \mathcal{L}_{\text{ranking}}(\theta \mid u, v) \implies \theta^t = \theta^{t-1} - \nabla_{\theta} \mathcal{L}_{\text{ranking}} = \theta^{t-1} - u_i \nabla_{\theta} \mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i) \quad (4)$$

Note how the gradients of the  $\mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i)$  wrt  $\theta$  are upweighted compared to vanilla SGD, by an exponentially decreasing factor of  $\mathcal{L}_{\text{CE}}(f_{\theta}(x_i), y_i)$ , i.e.,  $p_y^{1/\beta}$ . This helps the model to converge and rank datapoints with higher spuriousity before it moves on to those with lower levels of spuriousity.

**Ranking:** Finally, we set aside the samples for which  $v_i^{(t-1)*} = 1, v_i^{t*} = 0$ , and consider them as ranked by appending them to  $X_{\text{ranked}}^{t-1}$ , in decreasing order of their corresponding  $u_i$ s, thus allowing for fine-grained rank resolution. We provide the pseudocode for Sebra in Algorithm 1. [rebuttal: Based on the ordering obtained, we proceed in the next section, with formulating a contrastive learning based objective for learning a metric space devoid of spurious correlations.]

### 3.3 CONTRASTIVE DEBIASING

The ranking objective introduced in Section 3.1 generates a class-wise ordering of data points in the order of decreasing spuriousity. This fine-grained ranking can be leveraged for efficient debiasing. To demonstrate its effectiveness, we adopt a contrastive loss-based debiasing approach. [rebuttal: While contrastive learning has proven effective for debiasing (Zhang et al., 2022), the fine-grained bias characterization offered by Sebra enables the selection of more informative contrastive pairs. This approach surpasses traditional methods, which rely on simpler bias identification mechanisms, such as GCE or partially trained ERM models.] Additionally, contrastive learning-based approaches enable the utilization of the entire training dataset, unlike methods such as Kirichenko et al. (2023), which focus on the least spurious examples. Although the fine-grained bias identification generated by sebra could be integrated into other debiasing strategies like Kirichenko et al. (2023), we opt for a contrastive learning framework to showcase the full potential of Sebra.

Given a randomly sampled data point  $x_i$  with rank  $r$  from class  $c$ , we sample another instance  $x_n^-$  from the same class  $c$  and rank  $r$  to form a negative pair  $(x_i, x_n^-)$ . This is motivated by the ranking objective, which assigns the same rank to data points with similar levels of spurious correlations. To form a positive pair, we pair  $x_i$  with an instance of higher rank than  $r$ , as such instances are less likely to share the same spurious features as  $x_i$ . Using these contrastive pairs, we learn a debiased representation by optimizing the contrastive loss while simultaneously updating the full model via cross-entropy loss. For a classifier  $f_{\theta}$  with encoder  $f_{\text{enc}}$ , which maps a data point  $x$  to its representation  $z = f_{\text{enc}}(x)$ , the training objective is:

$$\hat{\mathcal{L}}(f_{\theta}; x, y) = \gamma \hat{\mathcal{L}}_{\text{con}}^{\text{sup}}(f_{\text{enc}}; x, y) + (1 - \gamma) \hat{\mathcal{L}}_{\text{CE}}(f_{\theta}(x, y), \quad (5)$$

where  $\gamma$  is a weighting coefficient. The supervised contrastive loss with  $M$  positive pairs and  $N$  negative pairs is:

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

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

## 4 EXPERIMENTAL SETUP

This section outlines the experimental framework for evaluating the effectiveness of the proposed spuriousity ranking and debiasing approach. We outline the datasets, evaluation metrics, and baselines used, with implementation details and hyperparameters in Appendix I.

**Datasets:** We evaluate the proposed ranking and debiasing strategy on one synthetic and two natural datasets with various spurious correlations. We use UrbanCars Li et al. (2023) for the synthetic dataset, focusing on car-type classification with spurious correlations involving the background and co-occurring objects. For natural datasets, we use CelebA Liu et al. (2015), which addresses spurious features like age and gender in predicting emotions (smiling / sad), following the setup in Hu et al. (2023). The presence of multiple spurious correlations and coarse-grained bias annotations in these datasets facilitates in defining a ground truth rank order of spuriousity as well as helps to evaluate the effectiveness of the proposed method in mitigating multiple biases. Furthermore, we test our method on one natural dataset: BAR (Nam et al., 2020), to demonstrate its scalability and effectiveness in natural settings. Sample images and dataset description are given in Appendix B.

**Evaluation Metrics :** Quantitatively comparing spuriousity rankings is challenging due to the absence of ground truth rankings. For datasets with bias annotations, such as Urban Cars and CelebA, where two biases are present (with Bias A being stronger than Bias B), we define the ground truth ordering as follows: (Bias A aligned, Bias B aligned), (Bias A aligned, Bias B conflicting), (Bias A conflicting, Bias B aligned), (Bias A conflicting, Bias B conflicting). We then compute Kendall’s tau correlation coefficient Kendall (1938) between this ground truth ordering and the rank orderings generated by our method and other baselines, quantitatively comparing ranking quality. In cases where even such coarse-grained bias annotations are unavailable as in BAR, we propose a quantitative metric, termed *Performance Disparity (PD)*. PD measures the difference in accuracy between models trained on the top-ranked images (highly spurious) and those trained on the bottom- $k$  ranked images (least-spurious), evaluated on an unbiased test set. Formally, PD is defined as:

$$PD = \text{Accuracy}_{\text{Bottom } k}(\mathcal{D}_{\text{test}}) - \text{Accuracy}_{\text{Top } k}(\mathcal{D}_{\text{test}}) \quad (6)$$

This metric captures the efficacy of the ranking scheme in segregating high and low spurious images to the head and tail of the output ranking. A high value of PD occurs when the bottom  $k$ -ranked data are the least spurious while the top-ranked data are highly spurious causing the models trained on these subsets to produce high and low test accuracy on an unbiased test set respectively. Thus, PD serves as a proxy to measure the quality of the obtained rankings. In addition to these quantitative assessments, we also present qualitative visualizations of the top-ranked and bottom-ranked images across all datasets in Appendix D, offering further insights into the nature of the spurious correlations and the effectiveness of the proposed ranking scheme.

To evaluate the robustness of the debiased model, we use a combination of conventional and new metrics. Li et al. (2023) introduce three novel metrics to evaluate debiasing in the presence of multiple spurious correlations in the UrbanCars dataset; BG Gap, CoObj Gap, and BG + CoObj Gap. These metrics are calculated based on the In-Distribution Accuracy (I.D.-Acc), representing the weighted average accuracy per group, with weights proportional to the frequencies of the groups in the training set. BG Gap, CoObj Gap, and BG + CoObj Gap measure the accuracy drop from ID-Acc to groups where the respective attributes are unaligned. Although initially proposed for the Urban Cars dataset, these metrics are versatile and can be applied to any dataset. For example, in CelebA, where the spurious correlations are Age and Gender, we calculate Age Gap, Gender Gap, and Age + Gender Gap. An average of these metrics (Avg. GAP) serves as an aggregate measure of robustness across different distribution shifts. In BAR, we report the test accuracy on the bias-conflicting set similar to prior methods (Li et al., 2022). We provide a detailed description of the metrics including their mathematical description in Appendix E.

**Baselines:** We compare the performance of the proposed approach with a supervised approach Group DRO (Sagawa et al., 2020) and five popular unsupervised approaches ERM (Vapnik, 1999), LfF (Nam et al., 2020), JTT (Liu et al., 2021), Debian (Li et al., 2022) and DFR (Kirichenko et al., 2023). The supervised approaches assume the availability of shortcut labels for all the spurious attributes while the unsupervised methods have access only to target labels. Both classes of methods further assume access to a small supervised validation set for hyperparameter tuning.

#### 4.1 RESULTS

In this section, we compare the performance of the proposed method with various baselines and datasets described in Section 4 to demonstrate its effectiveness in spuriousity ranking and debiasing.

**Ranking Evaluation.** As observed in Table 1, the proposed method produces a superior ordering of data points as indicated by a higher value of Kendall’s tau-b coefficient for both datasets. The

inferior performance of Spuriousity Ranking (Moayeri et al., 2023) despite human supervision could be attributed to the fact that biased attributes like background encompass multiple sub-attributes like lighting, sky, terrain, etc, and different sub-attributes are captured by different neurons rather than one or few neurons. Since these concepts are distributed across multiple neurons they need not be contained in top-k activations Fig. 4, resulting in many of these attributes being not considered while sorting the data, the rank ordering obtained via spuriousity ranking could be further improved by examining larger number of neurons at the expense of additional human supervision.

Table 2: Performance comparison across UrbanCars, CelebA, and BAR datasets. Sup.: Whether the model requires group or spurious attribute annotations (✗: not required, ✓: required). I.D. Acc. measures performance without subpopulation shift, while Avg GAP does so in its presence. All results are reported as mean (standard deviation).

Methods	Sup.	UrbanCars		CelebA		BAR
		I.D. Acc. (↑)	Avg GAP (↑)	I.D. Acc. (↑)	Avg GAP (↑)	Test Acc. (↑)
Group DRO	✓	91.60 (1.23)	-10.30 (1.35)	90.08 (0.70)	-5.79 (1.63)	-
ERM	✗	97.60 (0.86)	-31.90 (3.92)	96.43 (0.13)	-22.83 (0.84)	68.00 (0.43)
LfF	✗	97.20 (2.40)	-31.06 (3.56)	95.12 (0.35)	-22.57 (1.26)	68.30 (0.97)
JTT	✗	95.80 (1.45)	-20.50 (2.61)	91.86 (1.48)	-26.81 (2.53)	68.14 (0.28)
Debian	✗	98.00 (0.89)	-31.40 (1.44)	96.28 (0.37)	-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>-10.57 (1.72)</b>	88.61 (3.36)	<b>-9.82 (3.06)</b>	<b>75.36 (2.23)</b>

**Debiasing Evaluation.** As shown in Table 2, our proposed Sebra outperforms all the unsupervised methods in simultaneously mitigating multiple biases, as evidenced by the lower average gap (Avg. GAP) metric across datasets. While Sebra may not always achieve the best performance on individual Bias GAP metrics, this is due to the *whac-a-mole* dilemma observed in previous methods, where mitigating one bias attribute exceptionally well can amplify the other bias attribute. This can result in a very low Bias GAP for one attribute, even though the model remains highly biased overall. Furthermore, the proposed method consistently surpasses previous approaches. Additionally, our method performs comparably to prior single-bias unsupervised methods in single-bias settings, highlighting its effectiveness. An extended version of Table 2 including individual Bias GAP metrics is provided in Appendix D.2.

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 $\tau$ (↑)	Kendall’s $\tau$ (↑)	PD (↑)
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>35.47</b>

## 4.2 ANALYSIS AND ABLATION STUDIES

In this section, we present a comprehensive set of analyses and ablation studies to provide deeper insights into the performance of Sebra. Specifically, we investigate how the training dynamics of a model optimized using the proposed ranking objective differ from those of a standard empirical risk minimization (ERM)-based model. This comparison elucidates how the proposed selection and weighting mechanisms modulate the ERM training dynamics to facilitate spuriousity ranking. Furthermore, we conduct ablations on the various components of our framework to quantify their contributions to the overall ranking quality. Additional ablation studies are provided in Appendix F.

**Analysis of Ranking Dynamics:** In Section 3, we introduced Sebra, which integrates targeted modifications to ERM to systematically rank data points in the decreasing order of spuriousity. To rigorously assess the impact of these modifications, we conduct a detailed analysis of the training dynamics under the Sebra objective compared to standard ERM. Specifically, we leverage the UrbanCars dataset, which includes bias annotations, enabling a detailed evaluation of how spurious and intrinsic features are differentially learned across the two training paradigms. In Fig. 2, we plot the accuracy of three visual cues—object (e.g., car body type), background, and co-occurring objects—on the unbiased validation set by comparing the model’s {urban, country} predictions to the corresponding labels. As shown in Fig. 2 (Left), both models initially prioritize the easiest bias



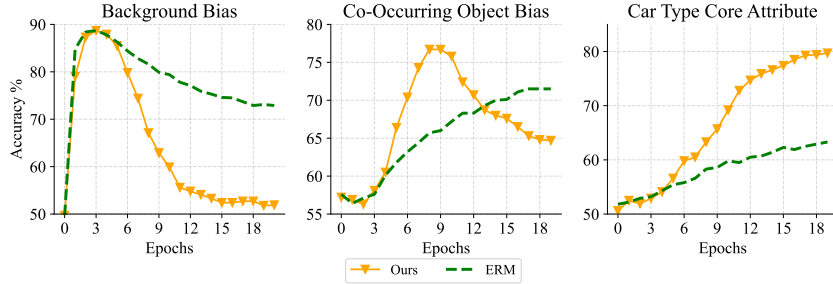


Figure 2: Training dynamics of Sebra and ERM monitored in terms of accuracies of background bias (left), co-occurring object bias (center), and core attribute (right).

attribute (background). However, as training progresses, the sebra objective induces a more pronounced forgetting of learned attributes compared to ERM, likely due to the selection mechanism through  $v_i$  that refocuses the model’s attention on other subgroups. The aggressive forgetting under the sebra framework overcomes the slowdown in the convergence of difficult attributes in the presence of simpler correlated features (Qiu et al., 2024), as evidenced by the higher peaks for both target and co-occurring object attributes. Another interesting observation is that, for relatively difficult bias attributes, such as co-occurring objects, the naive ERM formulation struggles to differentiate them from core attributes, as indicated by a simultaneous increase in accuracies in Fig. 2 (center and right). Sebra effectively addresses this challenge by leveraging the upweighting factor  $u_i$ , which amplifies the influence of highly spurious instances, thereby facilitating the progressive learning of these attributes. Additionally, the self-guided mechanism driven by  $u_i$  enhances overall performance, as demonstrated by the ranking improvements shown in Table 3. The non-overlapping peaks indicate that instances with a higher prevalence of respective attributes are assigned different ranks. Therefore, the ranking objective of sebra leads to well-segregated sequential learning of different shortcuts in the decreasing order of spuriousity. This analysis confirms our intuition and provides empirical evidence that the sebra efficiently ranks data in decreasing order of spuriousity.

**Effect of loss components:** To evaluate the contribution of various loss components, we conduct an ablation study by systematically removing components and measuring their impact on the quality of the resulting rankings using Kendall’s  $\tau$  coefficient. The results of this analysis are shown in Table 3. When using only  $\mathcal{L}_{CE}$ , corresponding to standard ERM training, we observe that ERM cannot alone rank data points effectively. To address this, we define a proxy ranking based on the epoch at which the predicted probability of the target attribute surpasses a fixed threshold. As shown in Table 3, the model trained with naive cross-entropy loss exhibits a low correlation with the ground truth ranking, as indicated by the low Kendall’s  $\tau$ . This suggests that naive cross-entropy fails to capture the underlying spuriousity of the data. The slightly positive  $\tau$  value may result from ERM’s inherent, though weak, ability to asynchronously learn different attributes. With the inclusion of  $v_i$  to the objective function, the ranking quality increases significantly, indicating that ERM’s poor bias ranking capability could be attributed to the interference caused by the easiest attributes in learning other attributes. Incorporating both  $v_i$  and  $u_i$  into the training objective improves ranking quality to 0.85, validating their importance.

Table 3: Ablation study of different components used in Sebra.

$\mathcal{L}_{CE}(\theta)$	$\mathcal{L}_{\text{ranking}}(v_i)$	$\mathcal{L}_{\text{ranking}}(u_i)$	Kendall’s $\tau$ ( $\uparrow$ )
✓	-	-	0.12
✓	✓	-	0.79
✓	✓	✓	<b>0.85 (Sebra)</b>

## 5 CONCLUSION AND FUTURE WORKS

We propose a novel debiasing strategy, *Sebra*, based on a fine-grained ranking of data points in decreasing order of spuriousity, obtained without any human supervision. Sebra facilitates spuriousity ranking by modulating the training dynamics of a simple ERM model to iteratively focus on highly spurious data points while simultaneously excluding already ranked datapoints from the ranking process. We further demonstrate how this fine-grained bias ordering enhances bias mitigation, by considering a contrastive learning-based approach as an exemplar on various datasets. Future work could explore bias mitigation strategies tailored to Sebra’s rankings, refine the ranking scheme, and develop unsupervised metrics for evaluating spuriousity rankings.

## REFERENCES

- C. Chang, G. Adam, and A. Goldenberg. Towards robust classification model by counterfactual and invariant data generation. In *CVPR*, 2021.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *ICLR*, 2019.
- Rui Hu, Yahan Tu, and Jitao Sang. Echoes: Unsupervised debiasing via pseudo-bias labeling in an echo chamber. In *Proceedings of the 31st ACM International Conference on Multimedia*. ACM, 2023.
- Minyoung Huh, Hossein Mobahi, Richard Zhang, Brian Cheung, Pulkit Agrawal, and Phillip Isola. The low-rank simplicity bias in deep networks. *TMLR*, 2023.
- Badr Youbi Idrissi, Martin Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data balancing achieves competitive worst-group-accuracy. In *CLear*, 2022.
- Yeonsung Jung, Hajin Shim, June Yong Yang, and Eunho Yang. Fighting fire with fire: Contrastive debiasing without bias-free data via generative bias-transformation. In *ICML*. PMLR, 2023.
- Maurice G Kendall. A new measure of rank correlation. *Biometrika*, 1938.
- Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to spurious correlations. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=Zb6c8A-Fghk>.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton Earnshaw, Imran Haque, Sara M Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. Wilds: A benchmark of in-the-wild distribution shifts. In *ICML*, 2021.
- Jungsoo Lee, Eungyeup Kim, Juyoung Lee, Jihyeon Lee, and Jaegul Choo. Learning representation via disentangled feature augmentation. In *NeurIPS*, 2021.
- Zhiheng Li, Anthony Hoogs, and Chenliang Xu. Discover and Mitigate Unknown Biases with Debiasing Alternate Networks. In *The European Conference on Computer Vision (ECCV)*, 2022.
- Zhiheng Li, Ivan Evtimov, Albert Gordo, Caner Hazirbas, Tal Hassner, Cristian Canton Ferrer, Chenliang Xu, and Mark Ibrahim. A whac-a-mole dilemma: Shortcuts come in multiples where mitigating one amplifies others. In *CVPR*, 2023.
- Yong Lin, Shengyu Zhu, Lu Tan, and Peng Cui. Zin: When and how to learn invariance without environment partition? In *NeurIPS*, 2022.
- Evan Z Liu, Behzad Haghighi, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *ICML*. PMLR, 2021.
- Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- Mazda Moayeri, Wenxiao Wang, Sahil Singla, and Soheil Feizi. Spuriousity rankings: Sorting data to measure and mitigate biases. In *NeurIPS*, 2023.
- Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure: Training debiased classifier from biased classifier. In *Advances in Neural Information Processing Systems*, 2020.
- Geon Yeong Park, Sangmin Lee, Sang Wan Lee, and Jong Chul Ye. Training debiased subnetworks with contrastive weight pruning. In *CVPR*, 2023.

- GuanWen Qiu, Da Kuang, and Surbhi Goel. Complexity matters: Feature learning in the presence of spurious correlations. In *ICML*, 2024.
- Shiori Sagawa, Pang Wei Koh\*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *ICLR*, 2020.
- Sahil Singla and Soheil Feizi. Salient imagenet: How to discover spurious features in deep learning? In *ICLR*, 2022.
- Nimit Sohoni, Jared Dunnmon, Geoffrey Angus, Albert Gu, and Christopher Ré. No subclass left behind: Fine-grained robustness in coarse-grained classification problems. In *NeurIPS*, 2020.
- Rishabh Tiwari and Pradeep Shenoy. Overcoming simplicity bias in deep networks using a feature sieve. In *ICML*, 2023.
- Vladimir Vapnik. An overview of statistical learning theory. *IEEE Trans. Neural Networks*, 1999.
- Yu Yang, Eric Gan, Gintare Karolina Dziugaite, and Baharan Mirzasoleiman. Identifying spurious biases early in training through the lens of simplicity bias. In *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*. PMLR, 2024.
- Yuzhe Yang, Haoran Zhang, Dina Katabi, and Marzyeh Ghassemi. Change is hard: A closer look at subpopulation shift. In *ICML*, 2023.
- Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Re. Correct-n-contrast: a contrastive approach for improving robustness to spurious correlations. In *ICML*, 2022.

## A PROOFS AND DERIVATIONS

### A.1 PROOF OF THEOREM 1

*Proof.* We start by first proving the case when  $u_i^* = e^{-t(\mathcal{L}_{\text{CE}}(x_i, y_i, \theta))}$  leads to the conservation law, i.e.,  $u_i^* = e^{-t(\mathcal{L}_{\text{CE}}(x_i, y_i, \theta))} \implies u_i \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) + \beta(u_i \ln u_i - u_i) = c$ .

$$\begin{aligned} \min_u u_i \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) - \lambda v_i^t + \beta g(u_i) &\implies \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) + \beta g'(u_i) = 0 \\ \implies g'(u_i) &= -\frac{1}{\beta} \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) \implies u_i^* = (g')^{-1}\left(-\frac{1}{\beta} \mathcal{L}_{\text{CE}}(x_i, y_i, \theta)\right) \end{aligned}$$

Now, we know that  $u_i^* = e^{-t(\mathcal{L}_{\text{CE}}(x_i, y_i, \theta))}$ . Considering  $t = \frac{1}{\beta}x$ , we have,  $(g')^{-1}(x) = e^x \implies g'(x) = \ln x$ . Then,

$$\begin{aligned} \ln u_i^* &= -\frac{1}{\beta} \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) \implies \int \ln u_i^* du_i = -\frac{1}{\beta} \int \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) du_i \\ &\implies u_i \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) + \beta(u_i \ln u_i - u_i) = c = \lambda v_i^t, \end{aligned}$$

since  $\lambda v_i^t$  was the constant that vanished under the derivative. This proves the statement in the direction  $u_i^* = e^{-t(\mathcal{L}_{\text{CE}}(x_i, y_i, \theta))} \implies u_i \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) + \beta(u_i \ln u_i - u_i) = c$ .

Next, we prove the statement in the other direction, i.e., when minimizing the conserved function leads to the exponentially decreasing characteristic of  $u_i^*$ . Solving for the minimizer of the conservation expression, we get:

$$\begin{aligned} u_i^* &= \min_u u_i \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) + \beta(u_i \ln u_i - u_i) - \lambda v_i^t \implies \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) + \beta \ln u_i = 0 \\ &\implies \ln u_i^* = -\frac{1}{\beta} \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) \implies u_i^* = e^{-\frac{1}{\beta} \mathcal{L}_{\text{CE}}(x_i, y_i, \theta)} = e^{-t(\mathcal{L}_{\text{CE}}(x_i, y_i, \theta))}, \end{aligned}$$

where  $t = \frac{1}{\beta}x$ . This proves the statement in the direction  $u_i \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) + \beta(u_i \ln u_i - u_i) = c \implies u_i^* = e^{-t(\mathcal{L}_{\text{CE}}(x_i, y_i, \theta))}$ , and completes the proof of the theorem.  $\square$

### A.2 SOLUTION FOR $v_i$

**Case 1 ( $k \geq 0$ ):** When  $k \geq 0$ , the optimal solution is  $v_i^{t*} = 1$ , ensuring  $\mathcal{L}(v | \theta, u) \geq 0$ . Otherwise,  $\mathcal{L}_{\text{ranking}}(v | \theta, u) = 0$ , which is always less than or equal to when  $v_i^t = 1$ . Thus,  $v_i^t = 1$  is the maximizer of  $\mathcal{L}_{\text{ranking}}(v | \theta, u)$  when  $k \geq 0$ .

Below, we derive the condition for optimality in terms of the predicted probability of the correct class  $p_y$  (this applies only when  $v_i^{t-1} = 1$ ):

$$u_i \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) - \lambda \geq 0 \implies u_i \ln p_y \leq -\lambda \implies p_y \leq e^{-\lambda/u_i}$$

**Case 2 ( $k < 0$ ):** When  $k < 0$ , the optimal solution is  $v_i^{t*} = 0$ . If  $v_i^t = 0$ ,  $\mathcal{L}_{\text{ranking}}(v | \theta, u) = v_i^t k$  would be negative. Thus,  $v_i^t = 0$  maximizes  $\mathcal{L}_{\text{ranking}}(v | \theta, u)$  in this case. Similarly to Case 1, we derive the condition for optimality when  $k < 0$  in terms of the predicted probability of the correct class  $p_y$ :

$$p_y > e^{-\lambda/u_i}.$$

### A.3 SOLUTION FOR $u_i$

$$\begin{aligned} \frac{\partial}{\partial u_i} \mathcal{L}_{\text{ranking}}(u | \theta, v) &= 0 \implies v_i^t \mathcal{L}_{\text{CE}}(x_i, y_i, \theta) - \beta + \beta[1 + \ln u_i] = 0 \\ \implies \ln p_y &= \beta \ln u_i \implies \ln u_i = \frac{1}{\beta} \ln p_y \implies \ln u_i = \ln p_y^{1/\beta} \implies u_i^* = p_y^{1/\beta} \end{aligned}$$

## B DATASETS

We evaluate the proposed method across three distinct datasets, each designed to explore different facets of bias and debiasing techniques. Below, we provide a succinct overview of each dataset:

1. **UrbanCars** Li et al. (2023): This synthetic dataset is purposefully crafted to investigate debiasing methodologies amidst multiple spurious correlations. Comprising two classes - UrbanCar and Country Car - each class encompasses 4000 samples. The dataset is characterized by two biased attributes: Background and Co-Occurring Object. UrbanCars feature city-like backgrounds with co-occurring objects such as traffic signs and fire hydrants, while Country Cars are set against rural backgrounds, predominantly featuring animals. UrbanCars is publicly available on Kaggle.
2. **CelebA**: A versatile dataset featuring celebrity faces alongside 40 binary attributes. We focus on the 'smile' attribute as the target, with biases introduced by age and gender. This configuration was introduced in Hu et al. (2023), and we employ their open-source code to obtain the data.
3. **Biased Action Recognition (BAR)** Nam et al. (2020): The Biased Action Recognition (BAR) dataset contains real-world images categorized into six action classes, each biased towards particular locations. The dataset includes six prevalent action-location pairs: Climbing on a Rock Wall, Diving underwater, Fishing on a Water Surface, Racing on a Paved Track, Throwing on a Playing Field, and Vaulting into the Sky. The testing set is composed exclusively of samples with conflicting biases. Therefore, achieving higher accuracy on this set signifies improved debiasing performance.

### B.1 DATA AUGMENTATIONS

For the CelebA dataset, we resize the images to a resolution of 224 x 224 and apply random horizontal flipping. For the Urban Cars dataset, we only apply random horizontal flip transformations. For BAR dataset, we apply both random horizontal flip and random resize crop. Note that these augmentations are applied during the contrastive debiasing stage while for the spuriousity ranking stage we apply no augmentations to prevent interference due to augmentations in effecting the spuriousity ranking.

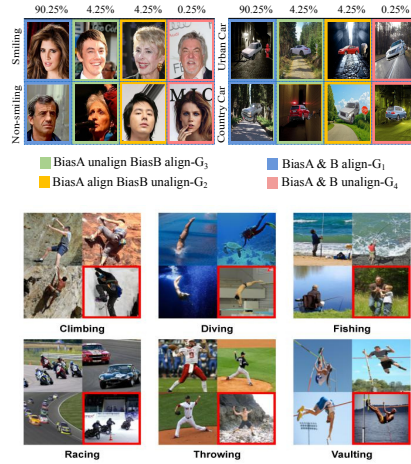


Figure 3: **Dataset samples:** Images from various datasets with multiple spurious correlations used in our experiments are shown below. For CelebA and UrbanCars dataset each column depicts multiple groups categorised based on biased features, as well as their proportions in the training set, each row displays samples from various classes. Images at the bottom demonstrates samples from BAR dataset from 6 classes. The images with red border lines belong to BAR evaluation set, and others belong to BAR training set.



## C BASELINES

We evaluate the proposed method against a series of unsupervised and supervised bias mitigation techniques. Below, we provide a concise overview of each method:

1. **GroupDRO**: Sagawa et al. (2020) A supervised bias mitigation technique leveraging group labels to identify and mitigate biases across various groups in the training data. The objective is to minimize the worst group accuracy across the identified groups.
2. **ERM** Vapnik (1999): Empirical Risk Minimization, employing cross-entropy loss and  $l_2$  regularization.
3. **Learning from Failure (LfF)**: Nam et al. (2020) This approach utilizes the Generalized Cross-Entropy (GCE) loss to derive a bias-only model. Subsequently, it learns a debiased model by reweighting the bias-conflicting points to learn a debiased model.
4. **JTT**: Liu et al. (2021) This method uses a ERM model trained for few epochs and identifies the misclassifications obtained by the model as bias conflicting samples and is upweighted for debiased learning.
5. **Debian**: Li et al. (2022) Introducing a novel bias identification scheme relying on the equal opportunity violation criteria, followed by bias mitigation strategies.
6. **DFR**: (Kirichenko et al., 2023) demonstrates that ERM model captures non-spurious attributes even when trained with biased training data and thus simple last layer retraining with unbiased data is sufficient for debiasing.

## D RESULTS

### D.1 QUALITATIVE RESULTS



Figure 4: Top 5 spurious concepts discovered using Spuriousity rankings introduced in Moayeri et al. (2023). As observed, the identified neurons capture only a subset of features corresponding to the spurious attribute 'background'; thus, ranking relying on top-k highly activating neurons would only rely on partial characteristics of spurious features.



Figure 5: Qualitative Analysis on UrbanCars Dataset: Examples of top-ranked and bottom-ranked samples as ranked by Sebra, showcasing a range of samples from both the classes.

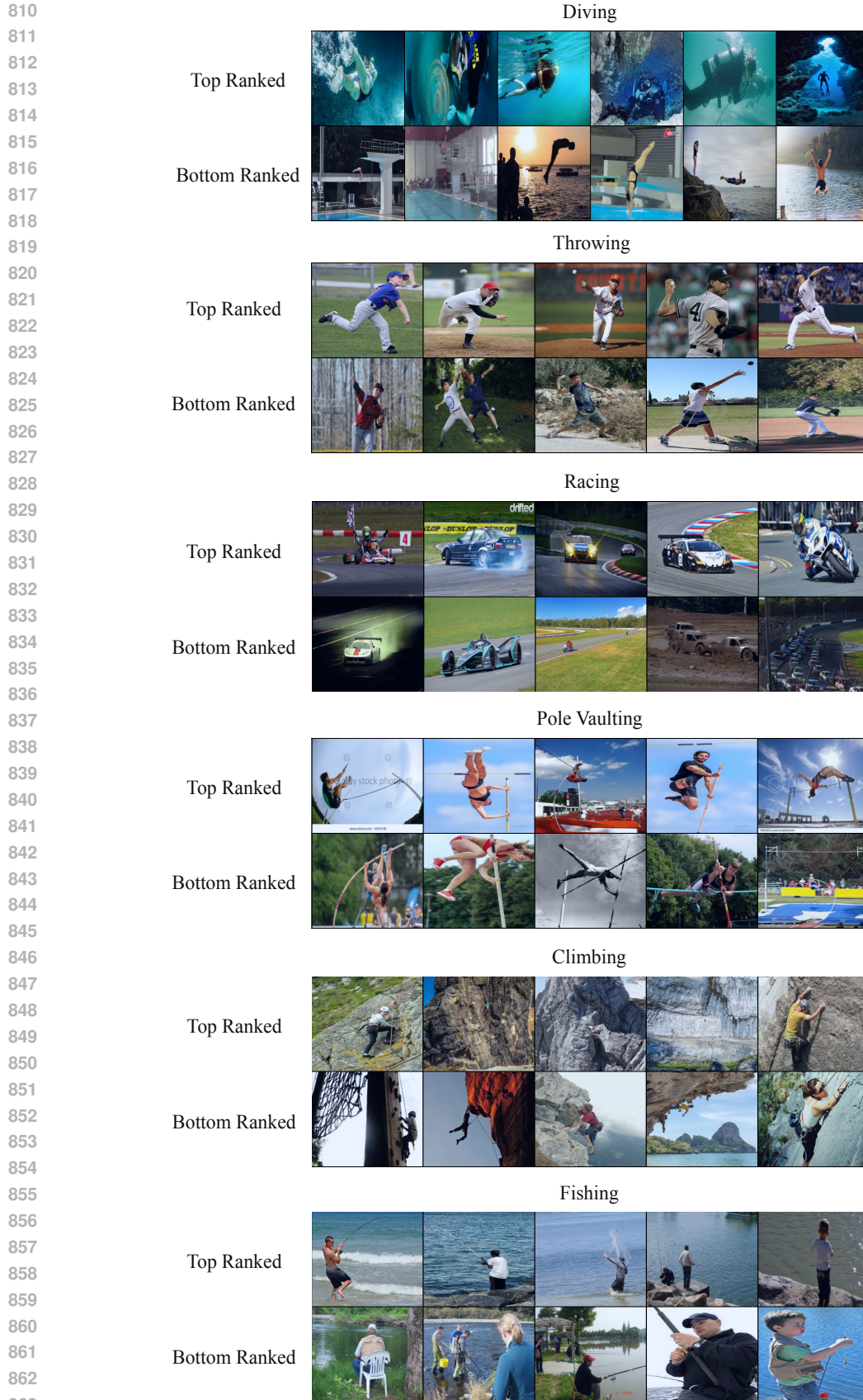


Figure 6: Qualitative Analysis on BAR Dataset: Examples of top-ranked and bottom-ranked samples as ranked by Sebra, showcasing a range of samples across different classes.

## D.2 QUANTITATIVE RESULTS

Table 4: Performance comparison on the UrbanCars dataset. Sup.: Whether the model requires group or spurious attribute annotations in advance (✗: not required, ✓: required). The best-performing results among unsupervised methods are marked in bold. The baseline results are taken from Li et al. (2023)

Methods	Sup.	I.D. Acc. (↑)	BG GAP (↑)	CoObj GAP (↑)	BG + CoObj GAP (↑)
Group DRO	✓	91.60(1.23)	-10.90 (1.08)	-3.60 (0.19)	-16.40 (2.80)
ERM	✗	97.60 (0.86)	-15.30 (1.35)	-11.20 (5.07)	-69.20 (5.34)
LfF	✗	97.20 (2.40)	-11.60 (1.23)	-18.40 (4.01)	-63.20 (2.21)
JTT	✗	95.80 (1.45)	-8.10 (1.08)	-13.30 (4.28)	-40.10 (2.48)
Debian	✗	98.00 (0.89)	-14.90 (1.08)	-10.50 (1.47)	-69.00 (1.78)
DFR	✗	89.70 (1.21)	-10.70 (1.85)	-6.90 (2.56)	-45.20 (3.42)
Sebra	✗	92.54 (2.10)	<b>-6.54 (1.38)</b>	-7.84 (1.38)	<b>-17.34 (2.40)</b>

Table 5: Performance comparison on the CelebA and BAR. Sup. indicates whether the method is supervised for bias (✓) or not (✗). The best results among unsupervised methods are marked in bold.

Methods	Sup.	CelebA				BAR
		I.D. Acc (↑)	Gender GAP (↑)	Age GAP (↑)	Gender+Age GAP (↑)	Test Acc. (↑)
Group DRO	✓	90.08 (0.70)	-5.67 (2.23)	-2.6 (2.4)	-9.11 (3.34)	-
ERM	✗	96.43 (0.13)	-22.7 (1.34)	-2.03 (0.77)	-43.77 (0.42)	68.00 (0.43)
LfF	✗	95.12 (0.35)	-24.14 (1.28)	-1.33 (1.2)	-42.26 (1.32)	68.30 (0.97)
JTT	✗	91.86 (1.48)	-31.07 (1.21)	-3.51 (2.44)	-45.85 (3.93)	68.14 (0.28)
Debian	✗	96.28 (0.37)	-22.03 (1.26)	-3.23 (1.65)	-42.41 (0.49)	69.88 (2.92)
DFR	✗	60.12 (1.28)	-12.16 (5.34)	-17.36 (3.23)	-27.96 (1.24)	69.22 (1.25)
Sebra	✗	88.61 (3.36)	<b>-2.21 (3.51)</b>	-6.89 (3.04)	<b>-20.36 (2.64)</b>	<b>75.36 (2.23)</b>

## D.3 DATASETS WITH OUTLIERS AND LABEL NOISE

Training datasets often contain samples from various origins and are labeled by annotators with differing expertise and background knowledge. Consequently, it is common for training sets to include outliers or mislabeled samples. Incorporating such corrupted instances into the training process for downstream tasks can adversely affect model performance, depending on the degree and prevalence of label noise.

The proposed ranking scheme offers a natural mechanism to address these issues. Specifically, in datasets with outliers and mislabeled instances, Sebra tends to assign the highest ranks to such corrupted samples. This property facilitates the identification and segregation of noisy data. For example, in the Living17 dataset, we empirically demonstrate this effect by showcasing samples with the highest and lowest ranks across several classes, as illustrated in Fig. 7.

This segregation enables an efficient filtration process, mitigating the negative impact of noisy data on subsequent training and enhancing model robustness.

## D.4 EMPIRICAL VALIDATION OF HARDNESS-SPURIOSITY SYMMETRY

Lin et al. (2022) *et al.* have theoretically demonstrated that unsupervised bias discovery is fundamentally impossible without the incorporation of additional inductive biases or meta-data. In this work, we leverage the concept of *Hardness-Spuriosity Symmetry* as an inductive bias to derive a continuous measure of spuriosity. This symmetry has been explored in prior studies, such as Nam et al. (2020); Qiu et al. (2024). Here, we refine and formalize this concept, proposing a method to quantitatively assess spuriosity.

To empirically validate this assumption, we present a plot of the training loss for samples with and without spurious correlations, generated by training a model using Empirical Risk Minimization (ERM) on the Urbancars dataset. As shown in Fig. 8, samples containing spurious correlations (i.e., bias attributes) exhibit a rapid decrease in loss, whereas non-spurious samples, which lack such shortcut attributes, show a much slower decline in loss. This discrepancy provides empirical support for our hypothesis that the difficulty of learning from a sample is inversely related to its spuriosity.



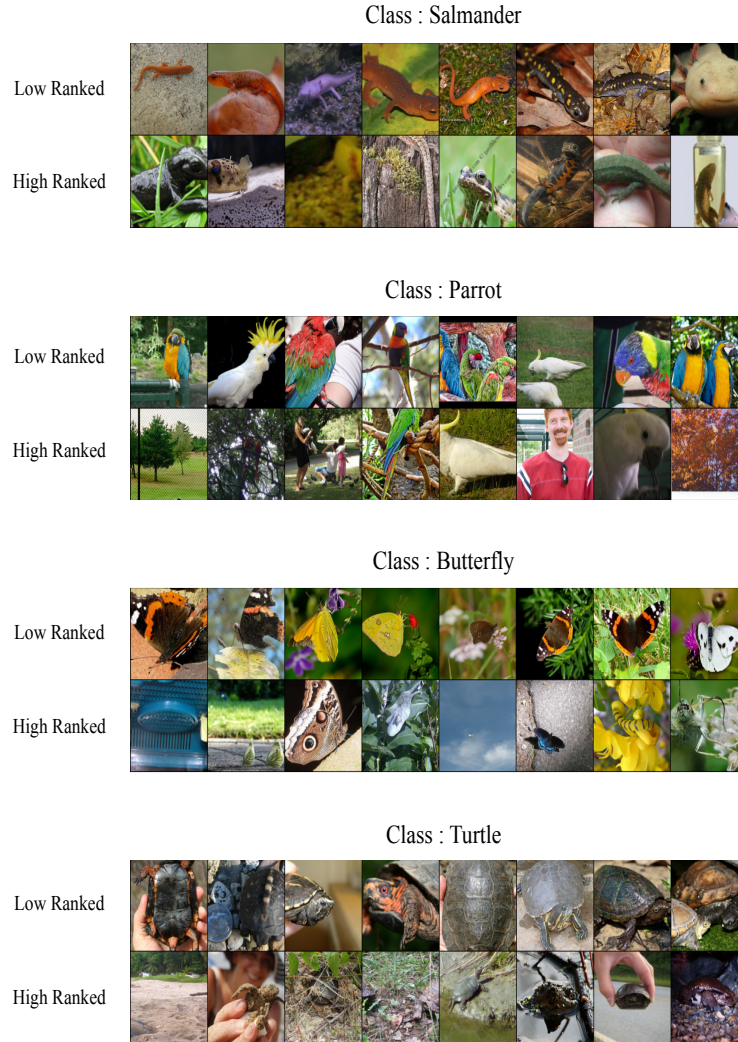


Figure 7: **Qualitative Analysis on the Living17 Dataset.** Examples of the least- and highest-ranked samples from select classes of the Living17 dataset. SEBRA assigns high ranks to mislabeled and outlier samples, facilitating their identification and removal during downstream task processing.

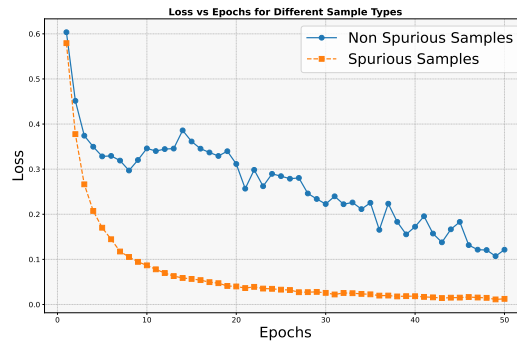


Figure 8: **Empirical Validation of Hardness-Spuriosity Symmetry:** Training loss vs. epochs for samples with and without spurious correlations on the Urbancars dataset. Samples with spurious correlations demonstrate a rapid decrease in loss compared to samples without such correlations, suggesting that higher spuriousity corresponds to easier learning.

## E METRICS

This section provides a detailed mathematical description of the evaluation metrics used throughout the paper.

1. **In-Domain Accuracy (I.D. Acc):** This metric represents the weighted average accuracy across groups, where the weights are determined by the correlation strength (i.e., frequency) of each group in the training data. It is designed to assess model performance under conditions where group distribution remains consistent with the training set.

$$\text{I.D. Acc} = \sum_{i=1}^G w_i \cdot \text{Acc}_i, \quad (7)$$

where  $w_i$  denotes the weight of group  $i$ , and  $\text{Acc}_i$  represents the accuracy for group  $i$ .

2. **Bias GAP:** This metric captures the difference between In-Domain Accuracy (I.D. Acc) and the accuracy on groups where the specific bias is less pronounced. It quantifies the model’s performance drop when tested on groups that diverge from the biases present in the training data.

$$\text{Bias GAP} = \text{I.D. Acc} - \text{Acc}_{\text{uncommon}}, \quad (8)$$

where  $\text{Acc}_{\text{uncommon}}$  represents the accuracy on groups with less prevalent bias.

3. **Kendall’s Tau Coefficient:** Kendall’s Tau is a non-parametric statistic that assesses the ordinal association between two variables. It ranges from -1 (perfect negative correlation) to 1 (perfect positive correlation), with 0 indicating no correlation. Particularly suitable for ranked data, Kendall’s Tau is more robust than Pearson’s correlation when the data distribution is non-normal or the relationship between variables is non-linear. The coefficient is computed by comparing the number of concordant and discordant pairs in the dataset.

## F ADDITIONAL ABLATION STUDIES

### F.1 EFFECT OF VARYING $\beta$

The Sebra objective introduced in Section 3.1 involves two key hyperparameters,  $\lambda$  and  $\beta$ . In this section, we investigate the sensitivity of the proposed ranking scheme to different values of  $\beta$ . Specifically, we plot the variation of Kendall’s  $\tau$  metric as a function of increasing  $\beta$  values in Fig. 9. As shown, the ranking quality demonstrates an almost linear decreasing trend as  $\beta$  increases, suggesting that smaller values of  $\beta$  are preferable for optimal performance. This behavior simplifies the hyperparameter search, as the optimal  $\beta$  appears to lie within the range (0, 1), reducing the computational cost associated with hyperparameter tuning.

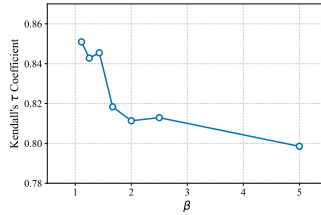


Figure 9: Sensitivity of ranking quality to  $\beta$

## G COMPUTATIONAL COST

The computational cost of debiasing via self-guided bias ranking can be divided into two components: the cost of spuriousity ranking with Sebra and the cost of contrastive debiasing. Sebra’s low computational complexity arises from the closed-form solution for the weighting variable and the progressive removal of data points during ranking, which accelerates the process. However, the cost of bias mitigation using contrastive learning is higher compared to simpler methods like JTT and Lff, due to its reliance on the full diversity of data rather than a limited subset. Despite this, Sebra enables fine-grained bias identification, making it compatible with various mitigation frameworks, and offers an efficient ranking process.



**Algorithm 1:** Pseudocode of Sebra

---

**Input:** A neural network  $f_\theta$ ,  $X_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N$ , where  $y_i \in \{1, \dots, C\}$ , maximum rank  $R$ , and upweighting and selection hyperparameters  $\beta$  and  $\lambda$ , respectively.

**Output:**  $X_{\text{ranked}} = \{(X_c, \rho(x_i))\}_{c=1}^C$

Initialize  $t = 0$

**while**  $t < R$  **do**

Obtain  $p_y = f_\theta(x, y)$  to compute  $u_i^*$  using equation 3  $\triangleright$ Up-weighting

Update the model parameters trained with upweighted points using equation 4  $\triangleright$ Training

Compute  $v_i^{t*}$  using equation 2 to select samples for subsequent training  $\triangleright$ Selection

**if**  $v_i^t = 0$  **and**  $v_i^{t-1} = 1$  **then**

$\rho(x_i) = t$   $\triangleright$ Ranking

Increment  $t = t + 1$

---

**H LIMITATIONS**

While Sebra demonstrates strong bias ranking capabilities and superior debiasing performance, it remains sensitive to label noise. Label noise can cause Sebra to incorrectly rank mislabeled samples as the least spurious, potentially compromising its effectiveness. Another limitation arises in datasets with multiple sub-population shifts, such as class imbalance. In such cases, the model may overemphasize a particular class, resulting in an increasingly unbalanced dataset during training. This imbalance can lead to learning collapse and a failure in ranking performance. Extending Sebra to handle these more complex scenarios, such as bias ranking in the presence of multiple sub-population shifts, could be a promising direction for future research.

**I REPRODUCIBILITY**

In this section, we outline the hyperparameters used in our proposed approach across various datasets. The optimal hyperparameters obtained for various datasets are summarised in Table 6. All experiments were conducted using a single RTX 3090 GPU. To facilitate reproducibility, we intend to release a user-friendly version of the code publicly along with the pre-trained models post-acceptance. We provide all implementation details and hyperparameters to facilitate reproducibility in Table 6. All the datasets used are publicly available or can be generated with publicly available resources.

**Implementation Details:** We use the same architectures and experimental setups as previous studies Li et al. (2022); Nam et al. (2020) to ensure fair comparisons. Specifically, we utilize ResNet-50 for the UrbanCars, and ResNet-18 for CelebA and BAR datasets. The optimal hyperparameters are selected based on experiments conducted on a small validation set with bias annotations, following the approach in Liu et al. (2021); Li et al. (2022) for CelebA and UrbanCars. For BAR, no bias annotations are used, even during validation and validation set is obtained by random split of training set in 80:20 ratio. To ensure statistical robustness, we perform four independent trials with different random seeds and report the mean and standard deviation of the results.

Table 6: Optimal hyper-parameters for the BAR , UrbanCars, and CelebA datasets determined through hyper-parameter search.

Parameter	UrbanCars	BAR	CelebA
Learning Rate (LR)	$1.0 \times 10^{-3}$	$1.0 \times 10^{-4}$	$1.0 \times 10^{-3}$
Batch Size	128	256	64
Optimiser	SGD	Adam	Adam
$\lambda$	0.75	0.9	0.8
$1/\beta$	0.7	0.8	0.8
Momentum	0.1	-	-
Weight decay	0.1	-	-