
Text-Guided Data Attribution: Attributing the Influence of Simplicity Bias to Dataset

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The effectiveness of deep learning models heavily relies on the quality and di-
2 versity of their training data. However, datasets collected from different sources
3 often introduce simplicity biases, where a models rely on easily learnable but
4 non-predictive (spurious) features for its predictions. While existing debiasing
5 techniques focus on model robustness, they leave the data untouched. Further,
6 they require manual group annotation of the entire training data or changes in
7 training strategy, which are often constrained by privacy, regulatory, or proprietary
8 constraints. As data becomes increasingly valuable, identifying and mitigating
9 bias directly at the data level has gained importance. Recently, data attribution
10 has emerged as a promising tool for uncovering issues in training data, yet its
11 vulnerability to simplicity bias has received limited attention. In this work, we
12 propose a novel data deletion framework that combines Neural Tangent Kernel
13 (NTK)-based data attribution with textual descriptions of bias to identify and re-
14 move training samples that do not significantly affect model performance. We first
15 demonstrate that NTK-based data attribution methods can themselves be influenced
16 by spurious features. Subsequently, to mitigate this, we use available metadata
17 or, when unavailable, a vision-language model, to annotate a small validation set
18 and extract a textual description of the bias. Based on this description, we identify
19 training samples that are semantically aligned with the spurious feature and exhibit
20 high detrimental attribution scores. Removing these samples from the training
21 data and retraining the model on the new training set improves its performance.
22 Our approach achieves better average and worst-group accuracy, outperforming
23 existing attribution-based baselines.

24 1 Introduction

25 The success of deep learning models is strongly influenced by the quality and quantity of the dataset
26 used for training [1–4]. These data are often collected via web scraping [5, 6], and external data
27 providers [7, 8]. However, such datasets can inadvertently contain illegal content [9] and can
28 encode negative societal biases [10, 11] that can influence model performance. In addition, data
29 collected from such varied sources can introduce distributional shifts, where subpopulations with
30 specific features may be overrepresented or underrepresented in the training data compared to the test
31 data [12].

32 These imbalances can introduce simplicity bias [13–15] where the model, due to high correlations
33 between specific features and the prediction task, relies on simpler, non-robust (spurious) features
34 instead of learning predictive features for classification. Several methods have been proposed to handle
35 such biases in the model. However, instead of addressing the data as the fundamental source of bias,
36 they primarily focus on improving model robustness by reweighting the loss function [16], modifying

the training objective [17–19], and model fine-tuning [20]. While effective, they inadvertently modify the standard training pipeline, which can increase the model’s susceptibility to adversarial attacks [21–24], and could conflict with regulatory requirements, especially in safety-critical settings [25, 26], which require adherence to a specific training regime for theoretical guarantees. Further, considering the inherent proprietary value of data [27] and the monetary investment needed for collecting a new dataset, it has become increasingly important to address these challenges at the data level.

A viable alternative in these scenarios could be to remove training samples containing spurious features [28, 29, 11], by ensuring that these samples don’t hurt the overall performance of the model, as in data attribution and Leave-One-Out (LOO) techniques [30–32]. Data attribution methods aim to estimate a model’s performance when specific training samples are excluded, enabling the evaluation of counterfactual scenarios—such as assessing the impact on test accuracy if certain subsets of the training data were omitted [31, 33, 34]. However, many of these methods are computationally expensive and can underperform in non-convex settings. Recent advancements in data attribution methods, such as Trak [32], leverage neural tangent kernels (NTK) to enable scalable data attribution for non-convex models [32]. However, the impact of spurious features on the data attribution scores generated by such methods remains an open question.

In our work, we demonstrate (Proposition 1, Appendix J) that in the presence of data bias, methods like Trak [32] can undervalue the attribution scores for training samples with spurious features [13–15]. This misattribution can hinder the identification of detrimental samples, especially for methods that rely solely on the magnitude of attribution scores [35, 31].

Motivated by these observations, we propose a two-stage strategy to mitigate the impact of spurious features - (a) In the first stage, we focus on identifying such features within the dataset using available meta-data or annotations generated by a vision language model. (b) In the second stage, we use multimodal embeddings, such as CLIP [36] to learn a metric [37, 38] that identifies training examples that are semantically similar to the spurious features identified in the first step and whose removal can improve the model’s performance as per the attribution scores.

The spurious features in the first stage are identified using metadata wherever available. In cases where metadata is unavailable, we utilize a vision-language model (VLM) to annotate a small validation set with its respective attributes and their associated values that are likely to introduce simplicity biases [39–41]. By evaluating the model’s performance on these attribute-value pairs and comparing it to the overall performance on the validation dataset [42], we identify potential spurious features and generate a corresponding textual description of these biases [43]. This textual representation enables targeted data pruning and helps to mitigate the impact of spurious features without relying on manual group annotations in the training dataset.

In summary, our contributions in this paper are as follows:

- We propose a novel data-centric approach that combines NTK-based data attribution methods with textual descriptions of underlying bias to mitigate the impact of spurious features in training datasets.
- We first theoretically demonstrate that NTK-based attribution scores can be influenced by spurious features, which may limit the effectiveness of methods that rely solely on these scores for data pruning. To overcome this limitation, we introduce a metric learning-based data deletion strategy that selectively removes training samples aligned with textual descriptions of spurious features and exhibiting low attribution scores.
- Our approach achieves up to a 4% gain in average accuracy, 18% in worst-group accuracy, and a 50% improvement in class-level performance across various datasets. Additionally, it outperforms NTK-based methods like Trak on average by 10.6% in worst-group accuracy for different biased datasets.

2 Related Work

2.1 Data Attribution

Data attribution methods provide a framework to relate a model’s predictions to its training dataset and have been used in a wide range of tasks, including model debugging and repair [44–47], subset selection [33, 34, 48], group robustness [11] and removing poisoning attacks [49].

The idea of linking a model’s predictions to its training data has been studied for decades under various names, including influence functions [50], regression analysis [51], and jackknife methods [52]. However, most of these early works focused on linear models and aimed to predict changes in the optimal parameters when individual or groups of samples were excluded during the learning process. Recent works have tried to extend influence function and jackknife-based attribution methods to non-linear models and bigger datasets [30, 53, 54]. However, despite their promising predictive capabilities, these methods often make strong assumptions of strong convexity and the existence of a unique global solution, which are not applicable for neural networks [55]. Furthermore, Basu et al. [56], Hammoudeh and Lowd [57] have demonstrated the fragile nature of methods like influence functions across different architectures, showing that they sometimes fail basic sanity checks. Various approaches have been proposed to address the limitations of influence functions, including gradient agreement scoring [58], training models to predict attribution scores, as in DataModels [59], and methods like Trak [32], which leverage concepts from the Neural Tangent Kernel (NTK) for data attribution. Unlike other approaches, such as DataModels, Trak does not require training thousands of models [32, 59] or tracking the loss changes over the entire training process, making it more efficient. However, the impact of spurious features within the dataset on the data attribution method like Trak remains largely unexplored.

2.2 Spurious Features and Simplicity Bias

Spurious features often arise from selection bias in the dataset [60], where, in the presence of multiple hypotheses for prediction, the model tends to rely on the simplest feature [61, 14, 13]. This preference can lead to suboptimal model performance, as it often ignores more robust and meaningful features that are essential for generalization in real-world scenarios. Various methods have been proposed to address spurious features in models. These include data augmentation techniques [62–67], and learning strategies that change the training objectives to make the model robust to spurious features [17, 68, 69, 19, 16, 70, 20]. However, many of these changes are restricted under the regulatory policy for safety-critical applications [71–74], especially considering privacy concerns associated with collecting datasets and model certification-based requirements [75, 76]. Recent work has explored data deletion as a strategy for mitigating spurious features [28, 29, 11]. These methods use group annotation of the dataset to remove random samples from majority groups [28, 29] or those with high detrimental attribution scores [11]. However, these methods often require manual group annotation of training [28, 29] or validation data [11], which is costly and time-consuming. Further, in real-world settings, where biases are identified post hoc after deployment and evolve over time [77], generating such annotations is often impractical, and enforcing a balance among different groups may result in excessive data removal from the majority group and can harm generalization [29]. Our method circumvents these limitations by using text-guided data attribution to efficiently remove harmful samples within a deletion budget, without relying on group labels or hurting model performance. Further details on limitations and capabilities of existing methods are discussed in Appendix E.

3 Proposed Method

3.1 Problem Definition

Consider a classification setting with a training dataset $\mathcal{D}_{\text{train}} = \{z_1, \dots, z_n\}$, where each sample $z_i = (x_i, y_i)$, consisting of an input (x_i) and associated class label (y_i) and a validation dataset, $\mathcal{D}_{\text{val}} = \{v_1, \dots, v_m\}$ with validation samples $v_j = (x_j, y_j)$. The training dataset ($\mathcal{D}_{\text{train}}$) is used to train a neural network with optimal parameters $\theta^*(\mathcal{D}_{\text{train}})$. Additionally, we assume that $|\mathcal{D}_{\text{val}}| \ll |\mathcal{D}_{\text{train}}|$.

Suppose for every training sample z there exists t underlying hidden discrete attributes, $A' = \{a^1, \dots, a^t\}$ and for each attribute (a^j) there are o possible values denoted as $V(a^j) \in \{b_1^j \dots b_o^j\}$. In real-world settings, neural networks (θ^*) trained on $\mathcal{D}_{\text{train}}$ often associate class labels (y) with specific attribute-value pairs (a^m, b_t^m) [43, 13, 14]. For example, a model trained to predict gender might associate it with the feature "beard" (present/absent). However, feature imbalance in the datasets can lead to misleading associations. If most of the male images in a dataset include smiles, the model might spuriously link "male" with "smiling" rather than "beard." This can cause misclassification,

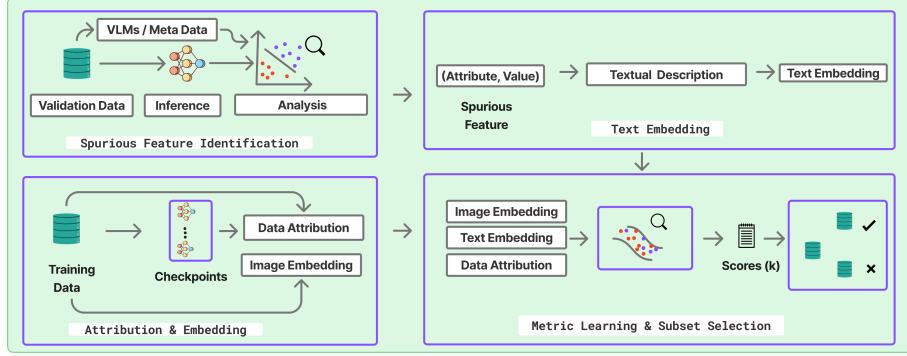


Figure 1: The figure illustrates the key steps in identifying detrimental samples. First, the performance of the model across different attribute value pairs is analyzed to identify and textually describe the underlying bias. Then, training samples that align with this bias and exhibit high detrimental attribution scores are selected for removal.

like predicting smiling females as males. We term such misleading attribute-value pairs as **spurious features**. In these scenarios, the primary objective of our work is to identify a set of detrimental examples, $\mathcal{S}^{\text{deter}} \in \mathcal{D}_{\text{train}}$, that “correspond” to the spurious features and might degrade the model’s performance. The model is then retrained from scratch on the filtered dataset, $(\mathcal{D}_{\text{train}} \setminus \mathcal{S}^{\text{deter}})$, to reduce the influence of the spurious feature in the training dataset similar to prior work like Chaudhuri et al. [28].

Our method for identifying $\mathcal{S}^{\text{deter}}$ involves two steps: (1) Annotate attribute–value pairs in the validation set to detect potential spurious features and generate a textual description of the bias; (2) Select $\mathcal{S}^{\text{deter}} \subset \mathcal{D}_{\text{train}}$ as samples semantically aligned with the bias and whose removal as per the data attribution scores does not degrade model performance.

3.2 Attribute Annotation and Spurious Feature Identification

A key component to identify spurious features is the availability of attribute–value annotations for the validation dataset. However, in many practical scenarios, such annotations are often missing from the metadata. Chen et al. [39] has shown that in the absence of such information, large language and vision models can be used to generate annotations necessary to identify the underlying spurious features. Hence, for datasets without pre-annotated attributes, we annotate the validation set with potential attribute–value pairs to assist in identifying spurious features.

To generate candidate attribute–value pairs, we leverage large language models such as ChatGPT [39]. ChatGPT is provided with a simple task description and prompted to suggest relevant attributes and associated values. For example, for a gender classification task, it can generate attributes like “smile”, “beard”, with possible values as “presence” or “absence”. We adopt task-specific prompts proposed by Chen et al. [39] to guide this process. Once the attribute–value pairs are generated, the next step is to annotate the validation dataset. However, considering the limitations associated with ChatGPT for this task [39], we use Llama 3.2 [78], a vision–language model, to annotate the images in the validation dataset. Further details about the prompts can be found in Appendix G.

3.2.1 Spurious Feature Identification

To identify spurious features, we take motivation from recent work that tries to identify systematic bias in a model [42, 43] based on its accuracy and errors on the validation dataset. However, unlike previous methods, which try to identify underperforming subgroups that may require collecting additional data, we try to determine the overperforming attribute-value pair as a possible candidate for data deletion [79, 28]. For this, we take inspiration from Johnson et al. [42] and compare the performance of the dataset associated with each attribute-value pair to the performance of the entire dataset. If the performance gap exceeds a predefined threshold, the corresponding attribute-value pair is flagged as a potential spurious feature in the model. Formally, this is expressed as:

$$\frac{1}{|\mathcal{D}_\alpha|} \sum_{(x,y) \in \mathcal{D}_\alpha} \mathbf{1}(h(x) = y) - \frac{1}{|\mathcal{D}_{\text{val}}|} \sum_{(x,y) \in \mathcal{D}_{\text{val}}} \mathbf{1}(h(x) = y) > \tau, \quad (1)$$

where, \mathcal{D}_α is a subset of validation data \mathcal{D}_{val} associated with a^v attribute and its j^{th} value b_j^v . The indicator function $\mathbf{1}$ indicates the correct prediction made by the model. The function $h(x)$ represents the prediction made by the model for a given input x , and y as the corresponding true class label. The parameter τ denotes the minimum threshold.

Once an attribute-value pair exceeds the threshold, a textual description is generated to describe the spurious feature. For example: "Images with $\{a\}$ as $\{b\}$." Here, (a, b) is the attribute value pair selected as per Equation 1. Details about the textual description are provided in Appendix H.4.

3.3 Coherent Data Attribution

After generating the desired text, the next task is to select a subset of data that is semantically coherent with the given text and whose removal can improve the performance of the model [80].

Since our task involves efficient subset selection, we formally define data attribution as follows:

Definition 1 (Data Attribution and Leave-one-out Influence Score [32]). *Given training dataset $\mathcal{D}_{\text{train}}$, and a model's utility function $\mathbf{f}(v; \theta)$ that measures the performance of the model, the data attribution score $\alpha : \mathcal{D}_{\text{train}} \times \mathcal{D}_{\text{val}} \rightarrow \mathcal{R}$ is defined as the change in the model's prediction for a validation sample v_i with respect to the optimal parameters when the training example z_k is excluded from the training dataset during the learning of the optimal parameters θ^* . Formally,*

$$\alpha(v_i; z_k) = \mathbf{f}\left(v_i; \theta^*(\mathcal{D}_{\text{train}})\right) - \mathbf{f}\left(v_i; \theta^*(\mathcal{D}_{\text{train}} \setminus z_k)\right) \quad (2)$$

For a classification task, the utility function $\mathbf{f}(z; \theta)$ for a sample $z = (x, y)$, [32], is defined as:

$$\mathbf{f}(z; \theta) = \log \left(\frac{p(z; \theta)}{1 - p(z; \theta)} \right), \quad (3)$$

where $p(z; \theta)$ represents the probability assigned to the correct class by the softmax function of a neural network parameterized by θ . A high $\mathbf{f}(z; \theta)$ corresponds to a high likelihood for a given sample (z) .

The NTK-based methods like Trak, have a closed-form formulation for data attribution score (α) (Definition 1) expressed as:

$$\begin{aligned} \alpha(v_j, z_i) &= \frac{1}{N} \sum_{n=1}^N \left(\phi_n(v_j)^\top (\Phi_n^\top \Phi_n)^{-1} \phi_n(z_i) \right) \\ &\times \frac{1}{N} \sum_{n=1}^N \left(1 - p_n^{z_i} \right) \end{aligned} \quad (4)$$

where, for N different checkpoints of model (θ^*), $\{p_n^{z_i}\}$ represents the probability assigned by n^{th} set of parameters to the correct class (y_i), for sample $z_i = (x_i, y_i)$. The terms $\phi_n(v_j)$ and $\phi_n(z_i)$ denote the projected gradients of the validation sample v_j and the training sample z_i with respect to the n^{th} set of optimal parameters, and for the utility function ($\mathbf{f}(\cdot; \theta_n^*)$). Additionally, Φ_n is the projected gradient for the entire training dataset. Further details about Trak can be found in Appendix B.

To quantify the impact of removing a data sample z from the training dataset on the performance of the entire validation dataset, we define the metric $\mathcal{A}(z)$ as a detrimental attribution score associated with the validation dataset for sample z . This metric measures the change in the model's performance (\mathbf{f}) for the validation dataset when z is excluded from the training dataset.

$$\begin{aligned} \mathcal{A}(z_i) &= - \sum_{v_j \in \mathcal{D}_{\text{val}}} \alpha(v_j, z_i) \\ &= \sum_{v_j \in \mathcal{D}_{\text{val}}} \left(\mathbf{f}(v_j; \theta^*(\mathcal{D}_{\text{train}} \setminus z_i)) - \mathbf{f}(v_j; \theta^*(\mathcal{D}_{\text{train}})) \right) \end{aligned} \quad (5)$$

where $z_i \in \mathcal{D}_{\text{train}}$. Unlike the data attribution score defined in Definition 1, $\mathcal{A}(z_i)$ is the negative of the general definition and evaluates the contribution of each training sample to the likelihood of the entire validation dataset. A higher value of $\mathcal{A}(z_i)$ indicates that removing the training sample z_i and retraining the model with the updated dataset leads to an optimal parameter θ^* that improves the likelihood of the validation dataset (Equation 3). In other words, training examples that degrade overall validation performance are assigned higher $\mathcal{A}(z_i)$ values. Once $\mathcal{A}(z_i)$ is calculated, it is normalized and used for further steps.

While removing samples with high $\mathcal{A}(z)$ values can improve the model’s performance; however, its impact on the downstream model is often tied up with its capability to remove samples with spurious features. During training, spurious features present in the dataset can result in gradient starvation [81, 61], a phenomenon that can hamper the learning of predictive features. Under such scenarios, we theoretically show that the detrimental attribution score (\mathcal{A}) for a data sample containing a spurious feature (f_1) can be lower than that of a data sample with predictive features (f_2), even when both features are equally represented. Consequently, deletion strategies based solely on high attribution scores may inadvertently remove examples with predictive rather than spurious features (Proposition 1) and can fail to capture the impact of removing data associated with spurious features on the overall generalization.

Proposition 1 (Under Valuation of Attribution Scores). *Consider a neural network in the neural tangent kernel (NTK) regime, trained using binary cross-entropy loss with two equally informative features, f_1 and f_2 . Let us assume that due to learning dynamics f_1 becomes dominant and causes gradient starvation of f_2 as per Pezeshki et al. [61]. Then, for two training samples z_i and z_j with equal representation of dominating features f_1 and f_2 , respectively. The attribution score for z_i can be systematically undervalued relative to z_j . Formally:*

$$|\mathcal{A}(z_i)| < |\mathcal{A}(z_j)|$$

The proof of Proposition 1, along with further details on gradient starvation, is provided in Appendix F. Empirical evidence supporting this phenomenon is presented in Appendix J.

This limitation of attribution scores motivates the need for a targeted removal strategy that specifically identifies and eliminates training samples sharing similar spurious features and exhibiting high $\mathcal{A}(z)$ scores. In many practical scenarios, the information about spurious features is missing in the data. Although annotating the entire training dataset using VLM-based models is possible, this approach is often excessively time-consuming and practically infeasible, particularly for large-scale datasets [40]. To address this, we adopt a zero-shot approach [82] and leverage textual descriptions of bias and CLIP embeddings to select data samples that are semantically similar to the identified textual descriptions. Specifically, we convert the textual description (Section 3.2) of the potential spurious feature into an embedding $\mathcal{C}_{\text{text}}$. Similarly, we convert all images in the training dataset into their corresponding CLIP embeddings $\mathcal{C}_{\text{image}}^i$ for $i \in 1, \dots, |\mathcal{D}_{\text{train}}|$. Each training sample z_i is then assigned a score k_i , reflecting its semantic similarity to the identified bias as per the given equation :

$$k_i = \exp \left(- \frac{(\mathcal{C}_{\text{text}} - \mathcal{C}_{\text{image}}^i) M (\mathcal{C}_{\text{text}} - \mathcal{C}_{\text{image}}^i)^\top}{2} \right),$$

where, $M = LL^\top, L \in \mathbb{R}^{D \times t}, t \ll D$ (6)

The text and image features, denoted as $\mathcal{C}_{\text{text}}, \mathcal{C}_{\text{image}}^i \in \mathbb{R}^{1 \times D}$ are represented as row vectors in a D-dimensional space. The matrix M is a positive semi-definite matrix, constructed as the outer product of a low-rank matrix L (rank at most t), and can serve as a learnable transformation. Since M defines the distance metric, varying the values of L allows us to generate different similarity measures for comparing data points [37, 38].

We aim to remove data samples that have high \mathcal{A} scores and are semantically aligned with the identified bias. To achieve this, we learn the matrix L [83, 37, 38] by maximizing the weighted \mathcal{A} score for each sample, where higher weights indicate stronger semantic alignment with bias as per Equation 6. To maintain semantic coherence with the bias description, the cumulative score for the dataset is enforced to exceed a threshold \mathcal{T} , defined as a fraction (β) of the total training size ($|\mathcal{D}_{\text{train}}|$). A larger \mathcal{T} emphasizes semantic alignment, while a smaller \mathcal{T} allows for flexibility in

sample selection based on \mathcal{A} scores. The complete optimization objective is described as below :

$$\begin{aligned} & \max_L \sum_{i=1}^{|\mathcal{D}_{\text{train}}|} \left(\frac{k_i}{\sum_j k_j} \right) \mathcal{A}(z_i) \\ \text{s.t. } & \sum_{i=1}^{|\mathcal{D}_{\text{train}}|} k_i \geq \mathcal{T}, \quad \mathcal{T} = \beta \times |\mathcal{D}_{\text{train}}|. \end{aligned} \quad (7)$$

To ensure that the optimization remains tractable, we replace the hard constraint with a soft penalty term [83] in the objective function. Further detail on this is provided in Appendix C.

Once the optimization is complete, a subset of training data with k_i scores greater than the hyperparameter γ is selected for removal ($\mathcal{S}^{\text{deter}}$). The model is then retrained with the updated training dataset ($\mathcal{D}_{\text{train}} \setminus \mathcal{S}^{\text{deter}}$) where, $\mathcal{S}^{\text{deter}} = \{z_i \in \mathcal{D}_{\text{train}} \mid k_i > \gamma\}$. A sensitivity analysis of all the hyperparameters, and comparison with only CLIP and only data attribution on overall performance is provided in Appendix Q and Appendix M, respectively.

4 Experiments

4.1 Setting

We evaluate the performance of our method across various datasets and compare it with existing data attribution techniques, including original training of model with complete dataset (original), Random deletion of data points (Random), Influence Function (IF) [30], TracIN [58], EWC Repair [31], and Trak [32]. The datasets used in our experiments include WaterBirds [84], Animal with attributes (AWA2) [85], German Traffic Sign Recognition Benchmark (GTSRB) [86], CELEBA [87, 43, 88] (Appendix H), CIFAR-10 [89], and ImageNet-100 [90, 91]. Further comparisons with robustness-based methods (groupDRO [17], JTT [16]) and group balancing methods are provided in Appendix I. For datasets such as GTSRB, CIFAR-10, and WaterBirds, we utilized attributes generated by ChatGPT and VLM models. To further assess the impact of metadata availability, we created two variants for the AWA2 datasets. The first variant, AWA2-A, includes class-specific annotations provided by the original datasets. The second variant, AWA2-B, uses attributes generated using ChatGPT and VLM-based annotation techniques (Section 3.2). All Primary experiments were conducted using a ResNet-18 model, which is the base architecture used in NTK-based data attribution methods such as Trak [32] for the image classification task. Additional experiments using alternative architectures and vision transformer models are presented in Appendix N and Appendix O, respectively. We have reported the worst group accuracy and average accuracy based on prior work on spurious features [67, 17, 28]. However, due to the absence of well-defined group structures in many real-world datasets [17], we have compared these datasets on average accuracy and class-level accuracy. All the experiments were conducted on two NVIDIA A6000 GPUs. Further details on training, hyperparameters, and subset size are provided in Appendix H. Algorithm 1 (Appendix) illustrates the overall workflow of our approach. We also report time and memory overheads associated with subset selection in Appendix R and Appendix S, respectively. Sample images from the selected subset $\mathcal{S}^{\text{deter}}$ are shown in Appendix U.

4.2 Improvement in Average Accuracy

Table 1 reports the improvement in average accuracy achieved by our method compared to existing baselines. On average, our method outperforms Trak by 1.4%, EWC by 1.6%, TracIN by 1.4%, Influence Functions by 2.0%, and the original full-dataset training baseline by 1.7%. Notably, we observe gains of 1.9%, 2.5%, and 2.4% over Trak on AWA2-B, WaterBirds, and AWA2-A, respectively. The performance improvement highlights the efficiency of our method in removing the detrimental samples associated with spurious features. We further saw a substantial improvement in under represented class as discussed in Section 4.3. Additional experiments on worst-group accuracy and architectural ablations for WaterBirds are provided in Appendix N.

4.3 Class Level Improvement after Data Deletion

Table 2 presents class-level accuracy for datasets with more than two classes. As per the results, our method improves the accuracy of a significant number of classes across datasets. For example, in

Table 1: Comparative evaluation of average accuracy of our proposed method (Ours) against baseline approaches across multiple datasets. The results report mean accuracy scores over three independent runs, with the best-performing values highlighted in **bold**. Entries with a gain of more than 1.5% over full-data training are highlighted in orange, while those exceeding 3% are shown in blue.

Dataset	Original	Random	IF	TracIN	EWC	Trak	Ours
WaterBirds	0.638	0.606	0.603	0.652	0.650	0.656	0.681
AWA2-A	0.644	0.622	0.644	0.652	0.642	0.638	0.662
CELEBA	0.895	0.893	0.890	0.893	0.890	0.898	0.906
GTSRB	0.969	0.966	0.973	0.971	0.975	0.971	0.980
AWA2-B	0.644	0.622	0.644	0.652	0.642	0.638	0.657
CIFAR-10	0.774	0.787	0.798	0.784	0.789	0.793	0.801
ImageNet-100	0.440	0.436	0.429	0.423	0.423	0.435	0.438

Table 2: Class-level accuracy improvement(Imp) after data removal across datasets. The table shows the maximum improvement in any class, the number of improved classes, and the mean improvement across them.

Dataset	Max Imp	# Imp Classes	Mean Imp
Awa2-A	16.27%	6 / 10	11.12%
Awa2-B	29.16%	4 / 10	17.98%
CIFAR-10	10.39%	7 / 10	5.59%
GTSRB	50.00%	22 / 43	5.69%
ImageNet-100	36.00 %	51 / 100	10.15%

Awa2-A, Awa2-B, ImageNet-100, and CIFAR-10, over 40% of the classes show improvement, with some achieving gains as high as 29.16%. Notably, in GTSRB, 22 out of 43 classes benefit, with a maximum per-class improvement of 50%. The improvement in average accuracy highlights that the improvement in underperforming classes is attained without substantially degrading the performance of other classes. Details on worst-class accuracy are provided in Appendix K.

Table 3: Comparison of best average accuracy across different data attribution methods for different spurious attributes. The table reports the mean accuracy across three independent runs. Entries with a gain of more than 1.5% over full-data training are highlighted in orange, while those exceeding 3% are shown in blue.

Target	Spurious Attribute	Original	Maj.-Rand	Random	IF	EWC	TracIN	Trak	Ours
arched eyebrows	receding hairline	0.713	0.740	0.739	0.716	0.724	0.730	0.722	0.736
attractive	mouth slightly open	0.628	0.627	0.668	0.640	0.633	0.631	0.658	0.673
big nose	male	0.771	0.770	0.770	0.764	0.751	0.745	0.756	0.780
goatee	bushy eyebrows	0.946	0.931	0.947	0.938	0.951	0.953	0.949	0.953
mouth slightly open	smiling	0.869	0.871	0.877	0.877	0.860	0.876	0.867	0.877
mouth slightly open	wearing lipstick	0.820	0.804	0.801	0.828	0.834	0.816	0.801	0.839
narrow eyes	eyeglasses	0.840	0.858	0.862	0.856	0.858	0.860	0.855	0.862
pointy nose	mouth slightly open	0.690	0.714	0.676	0.689	0.695	0.709	0.694	0.698
receding hairline	rosy cheeks	0.921	0.909	0.920	0.921	0.920	0.916	0.911	0.930
male	pointy nose	0.919	0.931	0.907	0.909	0.911	0.906	0.915	0.921

303

304 4.4 Performance across Different Spurious Attributes

To further investigate the impact of spurious features on both worst-group and average performance, we follow the setup of Eyuboglu et al. [43] and select a subset of the CELEBA dataset where the target attribute is strongly correlated with a spurious feature. We compare the average and worst-group performance achieved by our method against other baselines in Table 3 and Table 4. Additionally, considering the benefit of random data deletion in biased dataset [28] we introduce a new baseline, Maj.-Rand, where the subset of data is randomly deleted from the majority group. As shown in the results, our method outperforms other baselines in average accuracy in 7 and worst-group accuracy in 8 out of 10 settings, respectively. Notably, we observe a gain of over 4% in average accuracy for the target attribute attractive, compared to training on the original dataset. Similarly, worst-group accuracy improves by over 15% for attractive, receding hairline, and arched eyebrows, and by more than 5% for big nose, goatee, and male.

Table 4: Comparison of best worst-group accuracy across different data attribution methods for different spurious attributes. The table reports the mean accuracy across three independent runs. Entries with a gain of more than 5% over full-data training are highlighted in green, while those exceeding 15% are shown in violet.

Target	Spurious Attribute	Original	Maj.-Rand	Random	IF	EWC	TracIN	Trak	Ours
arched eyebrows	receding hairline	0.187	0.314	0.113	0.247	0.262	0.196	0.099	0.354
attractive	mouth slightly open	0.213	0.242	0.347	0.266	0.241	0.205	0.392	0.407
big nose	male	0.131	0.076	0.096	0.143	0.092	0.113	0.172	0.221
goatee	bushy eyebrows	0.432	0.493	0.287	0.437	0.439	0.387	0.278	0.548
mouth slightly open	smiling	0.524	0.415	0.552	0.418	0.441	0.487	0.433	0.489
mouth slightly open	wearing lipstick	0.555	0.471	0.557	0.598	0.594	0.549	0.486	0.612
narrow eyes	eyeglasses	0.208	0.052	0.119	0.000	0.092	0.128	0.024	0.151
pointy nose	mouth slightly open	0.045	0.044	0.046	0.034	0.028	0.021	0.040	0.084
receding hairline	rosy cheeks	0.121	0.228	0.131	0.179	0.241	0.254	0.201	0.296
male	pointy nose	0.840	0.882	0.824	0.833	0.861	0.870	0.875	0.903

Table 5: Comparison of best average and best worst-group accuracy between metadata-driven worst-group accuracy between our method and VLM-guided textual description.

Target	Sp. Attribute	Meta Data		VLM	
		Avg. Acc.	WG Acc.	Avg. Acc.	WG Acc.
bangs	black hair	0.922	0.649	0.916	0.624
big nose	wearing necklace	0.787	0.347	0.776	0.236
heavy makeup	straight hair	0.826	0.716	0.835	0.716
wearing earrings	bags under eyes	0.798	0.281	0.791	0.214

Table 6: Comparison of best average and best worst-group accuracy between our method and D3M across different spurious attributes.

Target	Spurious Attribute	Ours		D3M	
		Avg. Acc.	WG Acc.	Avg. Acc.	WG Acc.
bangs	black hair	0.922	0.649	0.920	0.627
big nose	wearing necklace	0.787	0.347	0.747	0.173
heavy makeup	straight hair	0.826	0.716	0.821	0.654
wearing earrings	bags under eyes	0.798	0.281	0.787	0.068

4.5 Ablation between Meta Data and VLM-based Description

Table 5 compares the performance of our method when using metadata versus VLM-generated textual descriptions of the spurious feature. While both strategies show comparable performance in terms of average accuracy, the metadata-driven variant generally achieves higher worst-group accuracy. This shows that a better annotation of underlying bias can help in the targeted removal of detrimental samples. However, even in the absence of such annotation, LLM and VLM-based methods can generate comparative performance.

4.6 Comparison with Group Annotation based Subset Selection

Table 6 presents a comparative evaluation between our method, which relies on the textual description of bias, against a technique that can use group annotation of spurious features in the validation dataset. To compare with such a method, we define group structure based on different values of Spurious Attribute and Target, and then use the method proposed by Jain et al. [11] (D3M) for subset selection. As per the result, on average, our method consistently outperforms D3M across both the best average and worst-group accuracy with a gain of 1.5% in best average accuracy and 11.8% in best worst group accuracy without using the explicit group annotation. This highlights the efficiency of the soft comparison scheme of clip features in handling partially visible features and the proposed optimization scheme compared to hard thresholding used in group annotation.

5 Conclusion

In this work, we propose a data deletion framework to mitigate the impact of spurious biases in the training dataset and enhance model performance. Our method employs metric learning techniques to target and remove training samples that are semantically aligned with the textual description of identified biases and whose removal, based on attribution scores, does not adversely affect model performance. To the best of our knowledge, this is the first approach to use text-guided data attribution scores to mitigate simplicity bias in models. However, its effectiveness depends on the quality of the textual descriptions used to capture spurious biases, and the current framework is limited to image datasets. In future work, we aim to incorporate a human-in-the-loop framework to better mitigate complex biases and to extend it to NLP tasks.

References

- [1] Nikita Bhatt, Nirav Bhatt, Purvi Prajapati, Vishal Sorathiya, Samah Alshathri, and Walid El-Shafai. A data-centric approach to improve performance of deep learning models. *Scientific Reports*, 14(1):22329, 2024.
- [2] Steven Euijong Whang, Yuji Roh, Hwanjun Song, and Jae-Gil Lee. Data collection and quality challenges in deep learning: A data-centric ai perspective. *The VLDB Journal*, 32(4):791–813, 2023.
- [3] Xinyi Xu, Shuaiqi Wang, Chuan-Sheng Foo, Bryan Kian Hsiang Low, and Giulia Fanti. Data distribution valuation. *arXiv preprint arXiv:2410.04386*, 2024.
- [4] Sang Keun Choe, Hwijee Ahn, Juhan Bae, Kewen Zhao, Minsoo Kang, Youngseog Chung, Adithya Pratapa, Willie Neiswanger, Emma Strubell, Teruko Mitamura, et al. What is your data worth to gpt? llm-scale data valuation with influence functions. *arXiv preprint arXiv:2405.13954*, 2024.
- [5] Zhipeng Xu, Zhenghao Liu, Yukun Yan, Zhiyuan Liu, Ge Yu, and Chenyan Xiong. Cleaner pretraining corpus curation with neural web scraping. *arXiv preprint arXiv:2402.14652*, 2024.
- [6] Jay M Patel and Jay M Patel. Introduction to common crawl datasets. *Getting structured data from the internet: running web crawlers/scrapers on a big data production scale*, pages 277–324, 2020.
- [7] Rodrigo F Berriel, Franco Schmidt Rossi, Alberto F de Souza, and Thiago Oliveira-Santos. Automatic large-scale data acquisition via crowdsourcing for crosswalk classification: A deep learning approach. *Computers & graphics*, 68:32–42, 2017.
- [8] Alexey Drutsa, Viktoriya Farafonova, Valentina Fedorova, Olga Megorskaya, Evfrosiniya Zermínova, and Olga Zhilinskaya. Practice of efficient data collection via crowdsourcing at large-scale. *arXiv preprint arXiv:1912.04444*, 2019.
- [9] David Thiel. Identifying and eliminating csam in generative ml training data and models. *Stanford Internet Observatory, Cyber Policy Center, December*, 23:3, 2023.
- [10] Emilio Ferrara. Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. *Sci*, 6(1):3, 2023.
- [11] Saachi Jain, Kimia Hamidieh, Kristian Georgiev, Andrew Ilyas, Marzyeh Ghassemi, and Aleksander Madry. Improving subgroup robustness via data selection. *Advances in Neural Information Processing Systems*, 37:94490–94511, 2024.
- [12] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International conference on machine learning*, pages 5637–5664. PMLR, 2021.
- [13] Rishabh Tiwari and Pradeep Shenoy. Overcoming simplicity bias in deep networks using a feature sieve. In *International Conference on Machine Learning*, pages 34330–34343. PMLR, 2023.
- [14] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.
- [15] RT McCoy. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. *arXiv preprint arXiv:1902.01007*, 2019.
- [16] 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 *International Conference on Machine Learning*, pages 6781–6792. PMLR, 2021.
- [17] Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019.
- [18] Jungsoo Lee, Eungyeup Kim, Juyoung Lee, Jiyeon Lee, and Jaegul Choo. Learning debiased representation via disentangled feature augmentation. *Advances in Neural Information Processing Systems*, 34: 25123–25133, 2021.
- [19] Sheng Liu, Xu Zhang, Nitesh Sekhar, Yue Wu, Prateek Singhal, and Carlos Fernandez-Granda. Avoiding spurious correlations via logit correction. *arXiv preprint arXiv:2212.01433*, 2022.

- 393 [20] Tyler LaBonte, Vidya Muthukumar, and Abhishek Kumar. Towards last-layer retraining for group
394 robustness with fewer annotations. *Advances in Neural Information Processing Systems*, 36:11552–11579,
395 2023.
- 396 [21] Gavin S Hartnett, Andrew J Lohn, and Alexander P Sedlack. Adversarial examples for cost-sensitive
397 classifiers. *arXiv preprint arXiv:1910.02095*, 2019.
- 398 [22] Kunyang Li, Jean-Charles Noiro Ferrand, Ryan Sheatsley, Blaine Hoak, Yohan Beugin, Eric Pauley, and
399 Patrick McDaniel. On the robustness tradeoff in fine-tuning. *arXiv preprint arXiv:2503.14836*, 2025.
- 400 [23] Pavan Kalyan Reddy Neerudu, Subba Reddy Oota, Mounika Marreddy, Venkateswara Rao Kagita,
401 and Manish Gupta. On robustness of finetuned transformer-based nlp models. *arXiv preprint*
402 *arXiv:2305.14453*, 2023.
- 403 [24] Chester Holtz, Tsui-Wei Weng, and Gal Mishne. Learning sample reweighting for accuracy and adversarial
404 robustness. *arXiv preprint arXiv:2210.11513*, 2022.
- 405 [25] Shiqi Wang, Kexin Pei, Justin Whitehouse, Junfeng Yang, and Suman Jana. Efficient formal safety
406 analysis of neural networks. *Advances in neural information processing systems*, 31, 2018.
- 407 [26] Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the convex outer
408 adversarial polytope. In *International conference on machine learning*, pages 5286–5295. PMLR, 2018.
- 409 [27] Feng Xiong, Maoyue Xie, Lingjuan Zhao, Cheng Li, and Xuan Fan. Recognition and evaluation of data
410 as intangible assets. *Sage Open*, 12(2):21582440221094600, 2022.
- 411 [28] Kamalika Chaudhuri, Kartik Ahuja, Martin Arjovsky, and David Lopez-Paz. Why does throwing away
412 data improve worst-group error? In *International Conference on Machine Learning*, pages 4144–4188.
413 PMLR, 2023.
- 414 [29] Badr Youbi Idrissi, Martin Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data balancing
415 achieves competitive worst-group-accuracy. In *Conference on Causal Learning and Reasoning*, pages
416 336–351. PMLR, 2022.
- 417 [30] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In
418 *International conference on machine learning*, pages 1885–1894. PMLR, 2017.
- 419 [31] Ryutaro Tanno, Melanie F Pradier, Aditya Nori, and Yingzhen Li. Repairing neural networks by leaving
420 the right past behind. *Advances in Neural Information Processing Systems*, 35:13132–13145, 2022.
- 421 [32] Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak:
422 Attributing model behavior at scale. *arXiv preprint arXiv:2303.14186*, 2023.
- 423 [33] Logan Engstrom, Axel Feldmann, and Aleksander Madry. Dsdm: Model-aware dataset selection with
424 datamodels. *arXiv preprint arXiv:2401.12926*, 2024.
- 425 [34] Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. Less: Selecting
426 influential data for targeted instruction tuning. *arXiv preprint arXiv:2402.04333*, 2024.
- 427 [35] Max Marion, Ahmet Üstün, Luiza Pozzobon, Alex Wang, Marzieh Fadaee, and Sara Hooker. When less
428 is more: Investigating data pruning for pretraining llms at scale. *arXiv preprint arXiv:2309.04564*, 2023.
- 429 [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
430 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from
431 natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR,
432 2021.
- 433 [37] Daryl Lim and Gert Lanckriet. Efficient learning of mahalanobis metrics for ranking. In *International*
434 *conference on machine learning*, pages 1980–1988. PMLR, 2014.
- 435 [38] Usha Bhalla, Alex Oesterling, Suraj Srinivas, Flavio P Calmon, and Himabindu Lakkaraju. Interpreting
436 clip with sparse linear concept embeddings (splice). *arXiv preprint arXiv:2402.10376*, 2024.
- 437 [39] Muxi Chen, Yu Li, and Qiang Xu. Hibus: on human-interpretable model debug. *Advances in Neural*
438 *Information Processing Systems*, 36, 2024.
- 439 [40] Haoming Lu and Feifei Zhong. Can vision-language models replace human annotators: A case study
440 with celeba dataset. *arXiv preprint arXiv:2410.09416*, 2024.

- [41] Zhen Tan, Dawei Li, Song Wang, Alimohammad Beigi, Bohan Jiang, Amrita Bhattacharjee, Mansoorreh Karami, Jundong Li, Lu Cheng, and Huan Liu. Large language models for data annotation and synthesis: A survey. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 930–957, 2024.
- [42] Nari Johnson, Ángel Alexander Cabrera, Gregory Plumb, and Ameet Talwalkar. Where does my model underperform? a human evaluation of slice discovery algorithms. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 11, pages 65–76, 2023.
- [43] Sabri Eyuboglu, Maya Varma, Khaled Saab, Jean-Benoit Delbrouck, Christopher Lee-Messer, Jared Dunnmon, James Zou, and Christopher Ré. Domino: Discovering systematic errors with cross-modal embeddings. *arXiv preprint arXiv:2203.14960*, 2022.
- [44] Chih-Kuan Yeh, Joon Kim, Ian En-Hsu Yen, and Pradeep K Ravikumar. Representer point selection for explaining deep neural networks. *Advances in neural information processing systems*, 31, 2018.
- [45] Siyi Tang, Amirata Ghorbani, Rikiya Yamashita, Sameer Rehman, Jared A Dunnmon, James Zou, and Daniel L Rubin. Data valuation for medical imaging using shapley value and application to a large-scale chest x-ray dataset. *Scientific reports*, 11(1):8366, 2021.
- [46] Harshay Shah, Sung Min Park, Andrew Ilyas, and Aleksander Madry. Modeldiff: A framework for comparing learning algorithms. In *International Conference on Machine Learning*, pages 30646–30688. PMLR, 2023.
- [47] Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al. Studying large language model generalization with influence functions. *arXiv preprint arXiv:2308.03296*, 2023.
- [48] Anshuman Chhabra, Peizhao Li, Prasant Mohapatra, and Hongfu Liu. " what data benefits my classifier?" enhancing model performance and interpretability through influence-based data selection. In *The Twelfth International Conference on Learning Representations*, 2024.
- [49] Chenwang Wu, Defu Lian, Yong Ge, Zhihao Zhu, and Enhong Chen. Influence-driven data poisoning for robust recommender systems. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10): 11915–11931, 2023.
- [50] Frank R Hampel. The influence curve and its role in robust estimation. *Journal of the american statistical association*, 69(346):383–393, 1974.
- [51] Daryl Pregibon. Logistic regression diagnostics. *The annals of statistics*, 9(4):705–724, 1981.
- [52] Rupert G Miller. The jackknife-a review. *Biometrika*, 61(1):1–15, 1974.
- [53] Kamiar Rahnama Rad and Arian Maleki. A scalable estimate of the extra-sample prediction error via approximate leave-one-out. *arXiv preprint arXiv:1801.10243*, 2018.
- [54] Ryan Giordano, William Stephenson, Runjing Liu, Michael Jordan, and Tamara Broderick. A swiss army infinitesimal jackknife. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1139–1147. PMLR, 2019.
- [55] Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger B Grosse. If influence functions are the answer, then what is the question? *Advances in Neural Information Processing Systems*, 35:17953–17967, 2022.
- [56] Samyadeep Basu, Philip Pope, and Soheil Feizi. Influence functions in deep learning are fragile. *arXiv preprint arXiv:2006.14651*, 2020.
- [57] Zayd Hammoudeh and Daniel Lowd. Identifying a training-set attack’s target using renormalized influence estimation. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages 1367–1381, 2022.
- [58] Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data influence by tracing gradient descent. *Advances in Neural Information Processing Systems*, 33:19920–19930, 2020.
- [59] Andrew Ilyas, Sung Min Park, Logan Engstrom, Guillaume Leclerc, and Aleksander Madry. Datamodels: Predicting predictions from training data. *arXiv preprint arXiv:2202.00622*, 2022.
- [60] Wenqian Ye, Guangtao Zheng, Xu Cao, Yunsheng Ma, and Aidong Zhang. Spurious correlations in machine learning: A survey. *arXiv preprint arXiv:2402.12715*, 2024.

- [61] Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. *Advances in Neural Information Processing Systems*, 34:1256–1272, 2021.
- [62] Megha Srivastava, Tatsunori Hashimoto, and Percy Liang. Robustness to spurious correlations via human annotations. In *International Conference on Machine Learning*, pages 9109–9119. PMLR, 2020.
- [63] Aahlad Puli, Nitish Joshi, Yoav Wald, He He, and Rajesh Ranganath. Nuisances via negativa: Adjusting for spurious correlations via data augmentation. *arXiv preprint arXiv:2210.01302*, 2022.
- [64] Huaxiu Yao, Yu Wang, Sai Li, Linjun Zhang, Weixin Liang, James Zou, and Chelsea Finn. Improving out-of-distribution robustness via selective augmentation. In *International Conference on Machine Learning*, pages 25407–25437. PMLR, 2022.
- [65] Xiangji Zeng, Yunliang Li, Yuchen Zhai, and Yin Zhang. Counterfactual generator: A weakly-supervised method for named entity recognition. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7270–7280, 2020.
- [66] Shirley Wu, Mert Yuksekgonul, Linjun Zhang, and James Zou. Discover and cure: Concept-aware mitigation of spurious correlation. In *International Conference on Machine Learning*, pages 37765–37786. PMLR, 2023.
- [67] Junhyun Nam, Jaehyung Kim, Jaeho Lee, and Jinwoo Shin. Spread spurious attribute: Improving worst-group accuracy with spurious attribute estimation. *arXiv preprint arXiv:2204.02070*, 2022.
- [68] Tan Wang, Chang Zhou, Qianru Sun, and Hanwang Zhang. Causal attention for unbiased visual recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3091–3100, 2021.
- [69] Daniel Levy, Yair Carmon, John C Duchi, and Aaron Sidford. Large-scale methods for distributionally robust optimization. *Advances in Neural Information Processing Systems*, 33:8847–8860, 2020.
- [70] Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to spurious correlations. *arXiv preprint arXiv:2204.02937*, 2022.
- [71] Eike Petersen, Yannik Potdevin, Esfandiar Mohammadi, Stephan Zidowitz, Sabrina Breyer, Dirk Nowotka, Sandra Henn, Ludwig Pechmann, Martin Leucker, Philipp Rostalski, et al. Responsible and regulatory conform machine learning for medicine: a survey of challenges and solutions. *IEEE Access*, 10:58375–58418, 2022.
- [72] Michael Matheny, S Thadaneey Israni, Mahnoor Ahmed, and Danielle Whicher. Artificial intelligence in health care: The hope, the hype, the promise, the peril. *Washington, DC: National Academy of Medicine*, 10, 2019.
- [73] Xudong Shen, Hannah Brown, Jiashu Tao, Martin Strobel, Yao Tong, Akshay Narayan, Harold Soh, and Finale Doshi-Velez. Towards regulatable ai systems: Technical gaps and policy opportunities. *arXiv preprint arXiv:2306.12609*, 2023.
- [74] Francisco Javier Campos Zabala. Responsible ai understanding the ethical and regulatory implications of ai. In *Grow Your Business with AI: A First Principles Approach for Scaling Artificial Intelligence in the Enterprise*, pages 453–477. Springer, 2023.
- [75] Romeo Valentin. Towards a framework for deep learning certification in safety-critical applications using inherently safe design and run-time error detection. *arXiv preprint arXiv:2403.14678*, 2024.
- [76] Ziquan Liu, Zhuo Zhi, Ilija Bogunovic, Carsten Gerner-Beuerle, and Miguel Rodrigues. Prosac: Provably safe certification for machine learning models under adversarial attacks. *arXiv preprint arXiv:2402.02629*, 2024.
- [77] Timothée Lesort. Spurious features in continual learning. In *AAAI Bridge Program on Continual Causality*, pages 59–62. PMLR, 2023.
- [78] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alan Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [79] Ming-Chang Chiu, Pin-Yu Chen, and Xuezhe Ma. Better may not be fairer: A study on subgroup discrepancy in image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4956–4966, 2023.

- [80] Long-Kai Huang, Peilin Zhao, Junzhou Huang, and Sinno Pan. Retaining beneficial information from detrimental data for neural network repair. *Advances in Neural Information Processing Systems*, 36, 2024.
- [81] Remi Tachet, Mohammad Pezeshki, Saeid Shabanian, Aaron Courville, and Yoshua Bengio. On the learning dynamics of deep neural networks. *arXiv preprint arXiv:1809.06848*, 2018. URL <https://arxiv.org/abs/1809.06848>.
- [82] Fei Pan, Sangryul Jeon, Brian Wang, Frank Mckenna, and Stella X Yu. Zero-shot building attribute extraction from large-scale vision and language models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 8647–8656, 2024.
- [83] Gabriel d’Eon, Jonathan d’Eon, James R. Wright, and Kevin Leyton-Brown. The spotlight: A general method for discovering systematic errors in deep learning models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1962–1981, New York, NY, USA, June 2022. ACM. doi: 10.1145/3531146.3534641.
- [84] Shiori Sagawa*, Pang Wei Koh*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=ryxGuJrFvS>.
- [85] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE transactions on pattern analysis and machine intelligence*, 41(9):2251–2265, 2018.
- [86] J. Stalkamp, M. Schlupsing, J. Salmen, and C. Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, 32:323–332, 2012. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2012.02.016>. URL <https://www.sciencedirect.com/science/article/pii/S0893608012000457>. Selected Papers from IJCNN 2011.
- [87] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [88] Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Ré. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations. *arXiv preprint arXiv:2203.01517*, 2022.
- [89] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [90] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- [91] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pages 776–794. Springer, 2020.
- [92] William B Johnson. Extensions of lipshitz mapping into hilbert space. In *Conference modern analysis and probability, 1984*, pages 189–206, 1984.
- [93] Zeyu Zhao. The application of the right to be forgotten in the machine learning context: From the perspective of european laws. *Cath. UJL & Tech*, 31:73, 2022.
- [94] Sungjin Lim and Junhyoung Oh. Navigating privacy: A global comparative analysis of data protection laws. *IET Information Security*, 2025(1):5536763, 2025.
- [95] Angela M Lonzetta and Thaier Hayajneh. Challenges of complying with data protection and privacy regulations. *EAI Endorsed Trans. Scalable Inf. Syst.*, 8(30):e4, 2021.
- [96] Oskar J Gstrein and Anne Beaulieu. How to protect privacy in a datafied society? a presentation of multiple legal and conceptual approaches. *Philosophy & Technology*, 35(1):3, 2022.
- [97] Jiaxin Fan, Qi Yan, Mohan Li, Guanqun Qu, and Yang Xiao. A survey on data poisoning attacks and defenses. In *2022 7th IEEE International Conference on Data Science in Cyberspace (DSC)*, pages 48–55. IEEE, 2022.
- [98] Ram Shankar Siva Kumar, David R O’Brien, Kendra Albert, and Salome Vilojen. Law and adversarial machine learning. *arXiv preprint arXiv:1810.10731*, 2018.

- 592 [99] Monty-Maximilian Zühlke and Daniel Kudenko. Adversarial robustness of neural networks from the
593 perspective of lipschitz calculus: A survey. *ACM Computing Surveys*, 57(6):1–41, 2025.
- 594 [100] Gilad Cohen, Guillermo Sapiro, and Raja Giryes. Detecting adversarial samples using influence functions
595 and nearest neighbors. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
596 recognition*, pages 14453–14462, 2020.
- 597 [101] Rabab Abdelfattah, Qing Guo, Xiaoguang Li, Xiaofeng Wang, and Song Wang. Cdul: Clip-driven
598 unsupervised learning for multi-label image classification. In *Proceedings of the IEEE/CVF international
599 conference on computer vision*, pages 1348–1357, 2023.
- 600 [102] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Conver-
601 gence and generalization in neural networks. In *Advances in Neural Information Process-
602 ing Systems*, volume 31, 2018. URL [https://proceedings.neurips.cc/paper/2018/file/
603 5a4be1fa34b8f9b53d01fdd7c1dc38c1-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/5a4be1fa34b8f9b53d01fdd7c1dc38c1-Paper.pdf).
- 604 [103] Greg Yang and Hadi Salman. A fine-grained spectral perspective on neural networks. *arXiv preprint
605 arXiv:1907.10599*, 2019.
- 606 [104] Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds,
607 Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqu Yan, et al. Captum: A unified and generic
608 model interpretability library for pytorch. *arXiv preprint arXiv:2009.07896*, 2020.
- 609 [105] Benedikt Boecking, Nicholas Roberts, Willie Neiswanger, Stefano Ermon, Frederic Sala, and Ar-
610 tur Dubrawski. Generative modeling helps weak supervision (and vice versa). *arXiv preprint
611 arXiv:2203.12023*, 2022.
- 612 [106] Olga Russakovsky and Li Fei-Fei. Attribute learning in large-scale datasets. In *European conference on
613 computer vision*, pages 1–14. Springer, 2010.
- 614 [107] Alan Pham, Eunice Chan, Vikranth Srivatsa, Dhruva Ghosh, Yaoqing Yang, Yaodong Yu, Ruiqi Zhong,
615 Joseph E Gonzalez, and Jacob Steinhardt. The effect of model size on worst-group generalization. *arXiv
616 preprint arXiv:2112.04094*, 2021.
- 617 [108] Muxi Chen, YU LI, and Qiang Xu. Hibug: On human-interpretable model debug. In *Thirty-seventh
618 Conference on Neural Information Processing Systems*, 2023. URL [https://openreview.net/forum?
619 id=4sDHLxKb1L](https://openreview.net/forum?id=4sDHLxKb1L).

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: The abstract and introduction clearly articulate the primary contributions of the paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: The paper explicitly discusses the limitations of existing methods in the Conclusion section, highlighting areas where current approaches fall short and talks about future research direction.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: All assumptions and detailed proofs are provided in Appendix F.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.

- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The codebase for reproducing the results is linked in the Abstract, and Appendix H provides additional implementation details.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The codebase for reproducing the results is linked in the Abstract, and Appendix H provides additional implementation details.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.

- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: The aforementioned details are available in Appendix H.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[Yes\]](#)

Justification: The mean across three runs is reported in the main draft.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [\[Yes\]](#)

Justification: Memory and time complexities are reported in Appendix S and Appendix R, respectively. Additional details can be found in the Experiments section.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.

- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Most ethical concerns are either not applicable to this work or have been appropriately addressed.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This work does not pose any specific societal risks; on the contrary, the proposed method actively addresses and mitigates model biases.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This work does not pose any such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.

- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Our code provides credit to external codebases in github repository.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We have not introduced any new assets in this work.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No crowdsourcing-based experiments were conducted in this work.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

899 Answer: [NA]
900 Justification: No experiments involving human subjects were conducted for this work
901 Guidelines:
902 • The answer NA means that the paper does not involve crowdsourcing nor research with human
903 subjects.
904 • Depending on the country in which research is conducted, IRB approval (or equivalent) may be
905 required for any human subjects research. If you obtained IRB approval, you should clearly state
906 this in the paper.
907 • We recognize that the procedures for this may vary significantly between institutions and
908 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for
909 their institution.
910 • For initial submissions, do not include any information that would break anonymity (if applica-
911 ble), such as the institution conducting the review.

912 **16. Declaration of LLM usage**

913 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard
914 component of the core methods in this research? Note that if the LLM is used only for writing,
915 editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or
916 originality of the research, declaration is not required.

917 Answer: [Yes]

918 Justification: Our work builds on the approach proposed by HiBug [39] and leverages large language
919 models (LLMs) and vision-language models (VLMs) for data annotation.

920 Guidelines:
921 • The answer NA means that the core method development in this research does not involve LLMs
922 as an important, original, or non-standard components.
923 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what
924 should or should not be described.

Algorithm 1 Proposed Method

Require: Training Dataset ($\mathcal{D}_{\text{train}}$), Validation Dataset ($\mathcal{D}_{\text{valid}}$), Number of Checkpoints (M), Rank for Metric Learning (t), Min Weight Fraction (β), Cutoff for Subset Selection (γ), CLIP Embedding Model (\mathcal{C}), Epochs for Classifier Training (\mathcal{E}), Optimization Iterations for Metric Learning (\mathcal{I}).

- 1: **## Classifier Training**
- 2: $i=0$
- 3: **for** epoch $\in [0 \dots \mathcal{E}]$ **do**
- 4: Train the classifier using $\mathcal{D}_{\text{train}}$.
- 5: **if** epoch $\in [\mathcal{E}, \mathcal{E} - 2, \mathcal{E} - 4, \mathcal{E} - 6, \mathcal{E} - 8]$ **then**
- 6: Save the checkpoint θ_i .
- 7: $i+=1$
- 8: **end if**
- 9: **end for**
- 10: Save $N = [\theta_0, \theta_1, \theta_2, \theta_3, \theta_4]$ checkpoints for the calculation of attribution score as per Equation 4.
- 11: **## Spurious Feature Identification**
- 12: Generate a list of possible attributes and corresponding values for $\mathcal{D}_{\text{valid}}$ using ChatGPT (Section 3.2).
- 13: **for** $i \in [1 \dots |\mathcal{D}_{\text{val}}|]$ **do**
- 14: Annotate attribute-value pairs for sample v_i using a Llama-based VLM model (Section 3.2).
- 15: **end for**
- 16: **## Calculating the Detrimental Attribution Score**
- 17: **for** $z_i \in \{z_1 \dots z_n\}$ **do**
- 18: Calculate the attribution score $\mathcal{A}(z_i)$ using the saved checkpoints (Equations 4 and 5).
- 19: **end for**
- 20: Compare the accuracy of each attribute-value pair using Equation 1. Flag an attribute-value pair as spurious if its accuracy exceeds the average dataset accuracy by a threshold τ .
- 21: Generate a textual representation of flagged attribute-value pairs under the context of the dataset (Appendix H.4).
- 22: Create a CLIP embedding of the textual representation ($\mathcal{C}_{\text{text}}$).
- 23: **## Metric Learning**
- 24: **for** $i \in [1 \dots |\mathcal{D}_{\text{train}}|]$ **do**
- 25: Calculate the CLIP image embedding ($\mathcal{C}_{\text{image}}^i$) for each sample z_i in $\mathcal{D}_{\text{train}}$.
- 26: **end for**
- 27: **for** $i \in [0 \dots \mathcal{I}]$ **do**
- 28: Optimize the loss \mathcal{L} using $\mathcal{C}_{\text{text}}$, $\mathcal{C}_{\text{image}}$, and $\mathcal{A}(z)$ as per Equation 10 to generate the metric \mathbf{k} using the hyperparameter t, β .
- 29: **end for**
- 30: Use the score \mathbf{k}, γ to identify $\mathcal{S}^{\text{deter}}$ and retrain the model on $\mathcal{D}_{\text{train}} \setminus \mathcal{S}^{\text{deter}}$.

926 **B Details on Trak**

$$\alpha(v_j, z_i) = \frac{1}{N} \sum_{n=1}^N \left(\phi_n(v_j)^\top (\Phi_n^\top \Phi_n)^{-1} \phi_n(z_i) \right) \times \frac{1}{N} \sum_{n=1}^N \left(1 - p_n^{z_i} \right)$$

where, $p_n^{z_i} = (1 + \exp(-y_i \mathbf{f}(x_i; \theta_n^*)))^{-1}$, $\phi_n(v_j) = \mathcal{P}^\top \nabla_{\theta} \mathbf{f}(v_j; \theta_n^*)$,
 $\phi_n(z_i) = \mathcal{P}^\top \nabla_{\theta} \mathbf{f}(z_i; \theta_n^*)$, $\Phi_n = [\phi_n(z_1)^\top; \dots; \phi_n(z_{|\mathcal{D}_{\text{train}}|})^\top]$
 $\Phi_n \in \mathbb{R}^{m \times k}$, $\mathcal{P} \sim \mathcal{N}(0, 1)^{p \times k}$, $k \ll p$. (8)

Equation 8, illustrates the calculation of the trak score. Scores consist of an average of the data attribution score calculated over multiple checkpoints (N). The terms $\phi_n(v_j)$ and $\phi_n(z_i)$ denote the projected gradients of the validation sample v_j and the training sample z_i for the n^{th} set of parameters and projection matrix \mathcal{P} . This projection matrix reduces the dimension of the gradient $\nabla_{\theta} \mathbf{f}(z; \theta_n^*) \in \mathbb{R}^p$ to a lower-dimensional space \mathbb{R}^k , where $k \ll p$, while approximately preserving the inner product, as per the classical Johnson-Lindenstrauss theorem [92].

C Soft Penalty for Optimization

For efficient optimisation of the constrained objective presented in Equation 7, we have replaced the hard constraint with a soft constraint $d(k)$ as per d'Eon et al. [83].

$$d(k) = C \cdot \max \left(\frac{\left(\sum_{i=1}^{|\mathcal{D}_{\text{train}}|} k_i - (\mathcal{T} + w) \right)^2}{w^2}, 0 \right), \quad (9)$$

This penalty term is quadratic and scaled by a shrinkable weight w , which is gradually reduced throughout the optimization process. The overall unconstrained optimization problem is defined in Equation 10 where C is a hyperparameter.

$$\mathcal{L} = \max_L \sum_{i=1}^{|\mathcal{D}_{\text{train}}|} \left(\frac{k_i}{\sum_j k_j} \right) \mathcal{A}(z_i) - d(k). \quad (10)$$

D Notations

Table 7: Notation table for key equation in main draft and proof

Symbol	Description
General Definitions	
$h(x)$	Model prediction for input x
y	True label corresponding to input x
$\mathbf{1}(h(x) = y)$	Indicator function: 1 if prediction is correct, else 0
\mathcal{D}_{val}	Validation dataset
$\mathcal{D}_{(a^v, b_j^v)}$	Subset of validation data with attribute-value pair (a^v, b_j^v)
$\mathcal{D}_{\text{train}}$	Training dataset: $\mathcal{D}_{\text{train}} = \{z_1, z_2, \dots, z_n\}$ where $z_i = (x_i, y_i)$
\mathcal{X}, \mathcal{Y}	Feature set $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$ and label set $\mathcal{Y} = \{y_1, y_2, \dots, y_n\}$
z_k, z_i	Training samples from $\mathcal{D}_{\text{train}}$
v_j	Validation sample
\hat{y}	Output of the final logit layer of a neural network
θ	Vectorized parameters of the neural network, $\theta \in \mathbb{R}^p$
e_i	Standard unit vectors
Neural Tangent Kernel (NTK) Specific	
$\mathcal{G}(\mathcal{X}, \theta)$	Neural Tangent Random Feature (NTRF) matrix: $\mathcal{G} = \frac{\partial \hat{y}(\mathcal{X}; \theta)}{\partial \theta}$
\mathcal{G}_0	NTRF matrix at initialization: $\mathcal{G}_0 = \mathcal{G}(\mathcal{X}, \theta_0)$
SVD Decomposition and Gradient Starvation	
U, S, V	Singular Value Decomposition (SVD) components: $Y \mathcal{G}_0 = U S V^\top$
u^i, v_k	Singular vectors from U and V corresponding to features
s_i	Singular value representing the strength of the i^{th} feature
Γ	Response of the network to features: $\Gamma = U^\top Y \hat{y} = S V^\top \theta$
Γ_i	Response of the i^{th} feature
Attribution and Trak Scoring	
Φ_n	Stacked gradient features of all training points for model n
$\mathbf{f}(v; \theta)$	Model output used for attribution (e.g., logit or loss) for input v under parameters θ
$\phi_n(\cdot)$	Projected gradient feature under model n
$\alpha(v_j, z_i)$	Attribution score: impact of removing z_i on prediction for v_j
$\mathcal{A}(z_i)$	Detrimental attribution score for training sample z_i
$p_n^{z_i}$	Predicted probability for z_i under model n
\mathcal{P}	Random projection matrix with entries drawn from $\mathcal{N}(0, 1)$
Optimization and others	
$C_{\text{text}}, C_{\text{image}}^i$	Text and image embeddings respectively
\mathcal{T}	Trade-off hyperparameter: $\mathcal{T} = \beta \times \mathcal{D}_{\text{train}} $ (control tradeoff between data attribution and semantic coherence)
β	Hyperparameter associated with \mathcal{T}
k_i	Selection weight for sample i
$d(k)$	Penalty term enforcing deletion constraint
\mathcal{L}	Final optimization objective including penalty
C	Hyperparameter associated with soft penalty Equation 9
$M = LL^\top$	Metric matrix constructed from L
$L \in \mathbb{R}^{D \times t}$	Learnable matrix under optimization defined by Equation 7
τ	Accuracy threshold to detect spurious bias
$\theta^*(\cdot)$	Final model parameters trained on the specified dataset

E Methods for Handling Spurious Features

Table 8 outlines the key capabilities and limitations of existing methods relative to ours. While data augmentation techniques [62–67] are widely adopted, they often require external data, which can conflict with privacy and regulatory constraints [93–96]. Moreover, without appropriate supervision, they risk introducing new spurious features or being vulnerable to data poisoning attacks [97, 98].

In contrast, methods such as group annotation-based optimization (e.g., gDRO [17]), loss reweighting techniques (e.g., JTT [16]), and final-layer fine-tuning [70, 20] do not pose privacy risks. However, in safety-critical applications where models must satisfy stability guarantees [75, 76, 99], these methods can compromise robustness, especially when models are required to ensure Lipschitz continuity for certification. Specifically, they are susceptible to targeted attacks [21–24], particularly when the training procedure heavily relies on a small subset of influential examples [100, 98] used for fine-tuning or reweighting loss values.

Group-balancing techniques [28, 29] partially address these challenges, but often over-prune majority groups. In contrast, our method supports budget-constrained, targeted sample removal, ensuring only detrimental examples are excluded during training.

Furthermore, many of these methods [17, 28, 29] rely on manual group annotations of the training dataset. As spurious features [77] evolve post-deployment, maintaining robustness would require repeated manual annotation cycles. In contrast, our approach eliminates the need for group labels for the training dataset and leverages textual descriptions of bias to guide targeted data removal. The use of a textual description of the bias and the proposed metric learning approach provides a zero-shot approach [101, 82] to approximate the underlying group structure without having any annotation overhead. This design also allows integration of feedback from subject-matter experts, making the process more adaptive and practical.

Table 8: Comparison of methods across regulatory and robustness capabilities.

Method	Regulatory Restrictions	Supports Textual Descriptions	No Group Annotation in Training Data	Privacy	Prevent Over pruning of Majority Group	Robust to Adv-Attacks
Data Augmentation	✗	✗	✓	✗	✓	✗
Group-Annotation based Optimization	✗	✗	✗	✓	✓	-
Reweighting Loss/Data	✗	✗	✓	✓	✓	✗
Last Layer Fine-Tuning	✗	✗	✓	✓	✓	✗
Group Balancing Method	✓	✗	✗	✓	✗	✓
Ours	✓	✓	✓	✓	✓	✓

F Theoretical Formulation

For dataset $\mathcal{D}_{\text{train}} = \{z_1 \dots z_n\}$ where $z_i = (x_i, y_i)$, and $x_i \in \mathbb{R}^d$, and corresponding labels $y_i \in \{-1, +1\}^n$. Let \hat{y} denotes the output of the final logit layer of an L-layer neural network trained using binary cross-entropy, and $\theta \in \mathbb{R}^p$ represents a p-dimensional vectorized parameter of the neural network (Equation 4, Equation 8). let $\mathcal{X} = \{x_1 \dots x_n\}$ and $\mathcal{Y} = \{y_1 \dots y_n\}$ constitute the respective features and class labels.

In the Neural Tangent Kernel (NTK) framework [102], the final output of a neural network can be approximated as a linear function of parameters, whose properties are governed by the Neural Tangent Random Feature (NTRF) matrix, defined as:

$$\mathcal{G}(\mathcal{X}, \theta) = \frac{\partial \hat{y}(\mathcal{X}; \theta)}{\partial \theta}, \quad \mathcal{G} \in \mathbb{R}^{n \times p}. \quad (11)$$

For wide-width neural networks, the NTRF matrix remains approximately constant during training [61], allowing the output of the neural network to be approximated using the initial NTRF matrix, $\mathcal{G}_0 = \mathcal{G}(\mathcal{X}, \theta_0)$, as follows:

$$\hat{y}(\mathcal{X}, \theta) = \mathcal{G}_0 \theta. \quad (12)$$

The dominant features of the dataset can be estimated using the principal components of $\mathcal{G}_0 = \mathcal{G}(\mathcal{X}, \theta_0)$, which are equivalent to the principal components of the NTK gram matrix [103].

Definition 2 (Features and gradient starvation [61]). *Consider a support vector decomposition of $Y\mathcal{G}_0 = USV^\top$, where $Y = \text{diag}(y)$, the i^{th} feature is represented by $(V^\top)_{(i,:)}$ or $(V)_{(:,i)}$ with its strength denoted as $s_i = (S)_{ii}$ and its weight across all training samples represented by $(U)_{(:,i)}$. The response of the neural network to the i^{th} feature can be expressed as Γ_i , where:*

$$\Gamma := U^\top Y \hat{y} = SV^\top \theta.$$

Due to the imbalance in the training dataset, for a given set of features and the optimal parameter θ^* , the presence of the i^{th} feature can influence the learning of the j^{th} feature. This phenomenon, referred to as gradient starvation, arises in optimal parameters if:

$$\frac{d\Gamma_j^*}{d(s_i^2)} < 0$$

Definition 2 suggests that as the strength of the i^{th} feature (s_i^2) increases, the learning of the j^{th} feature gets impacted. This implies that stronger features can dominate the learning process, leading to a reduced contribution of other informative features in the model's predictions.

Theorem 1 (Gradient Starvation Regime [61]). *For a neural network in the linear regime and trained using binary cross entropy loss with feature coupling between two features f_1 and f_2 as defined in Pezeshki et al. [61] and with $s_1^2 > s_2^2$, we have,*

$$\frac{d\Gamma_2^*}{d(s_1^2)} < 0,$$

Now, under the given setting, we will try to understand the influence of gradient starvation on the performance of the NTK-based data attribution methods :

Proposition 2 (Under Valuation of Trak Scores). *Consider a neural network in the neural tangent kernel (NTK) regime, trained using binary cross-entropy loss with two equally informative features, f_1 and f_2 . Let's assume that due to learning dynamics f_1 becomes dominant and cause gradient starvation of f_2 as per Pezeshki et al. [61]. Then, for two training samples z_i and z_j with equal representation of dominating features f_1 and f_2 respectively. The attribution score for z_i can be systematically undervalued relative to z_j . Formally:*

$$|\mathcal{A}(z_i)| < |\mathcal{A}(z_j)|$$

Proof. For a sigmoid-based activation, the output probability for feature set (\mathcal{X}) is given by:

$$\begin{aligned} p(\mathcal{X}; \theta) &= \frac{1}{1 + \exp(-\hat{y}(\mathcal{X}; \theta))}, \\ p(\mathcal{X}; \theta) \cdot (1 + \exp(-\hat{y}(\mathcal{X}; \theta))) &= 1, \\ p(\mathcal{X}; \theta) \cdot \exp(-\hat{y}(\mathcal{X}; \theta)) &= 1 - p(\mathcal{X}; \theta), \\ \hat{y}(\mathcal{X}; \theta) &= \log\left(\frac{p(\mathcal{X}; \theta)}{1 - p(\mathcal{X}; \theta)}\right). \end{aligned} \tag{13}$$

Hence, the utility function (f) used in Trak for data attribution (Equation 3) is equivalent to the logit of a binary cross entropy (\hat{y}).

From the definition of gradients:

$$\frac{\partial \hat{y}(x; \theta)}{\partial \theta} = \frac{\partial \mathcal{G}_0 \cdot \theta}{\partial \theta} = \mathcal{G}_0 \tag{14}$$

As per Equation 4 and Equation 8, $\Phi_m = \mathcal{G}_0 \cdot \mathcal{P}$. Now, considering that the projection matrix [92] preserves the inner product of the actual gradient vector. We will simplify our argument and calculate the value for the unprojected gradients [32] ($\mathcal{P} = I_d$). Furthermore, under the NTK regime, where the optimal parameters are similar [102], we calculate the attribution score for a single checkpoint ($M=1$). For ease of derivation, we will omit the subscript m i.e., $\Phi_1 = \Phi$ and $\phi_1 = \phi$, hence:

$$\Phi = \mathcal{G}_0 \tag{15}$$

$$\Phi^T \Phi = \mathcal{G}_0^T \mathcal{G}_0 \tag{16}$$

Now, as per the feature decomposition defined in Definition 2 :

$$\begin{aligned} Y \mathcal{G}_0 &= U S V^T \\ (Y \mathcal{G}_0)^T (Y \mathcal{G}_0) &= \left(U S V^T \right)^T \left(U S V^T \right) \\ \mathcal{G}_0^T Y^T Y \mathcal{G}_0 &= V S^2 V^T \end{aligned} \tag{17}$$

Since $Y = \text{diag}\{y_1, \dots, y_n\}$ and $y \in \{-1, 1\}$, it follows that:

$$\begin{aligned} Y^T Y &= I, \\ \mathcal{G}_0^T \mathcal{G}_0 &= V S^2 V^T, \\ \Phi^T \Phi &= V S^2 V^T. \end{aligned} \tag{18}$$

1004 The validation attribution score (Equation 5) is given by :

$$\begin{aligned}\mathcal{A}(z_i) &= \sum_{v_j \in \mathcal{D}_{val}} -\alpha(v_j; z_i) \\ &= \sum_{v_j \in \mathcal{D}_{val}} -\phi(v_j)^\top (\Phi^\top \Phi)^{-1} \phi(z_i) (1 - p^{z_i})\end{aligned}\tag{19}$$

1005 Substituting the value of $\Phi^\top \Phi$:

$$\begin{aligned}\mathcal{A}(z_i) &= \sum_{v_j \in \mathcal{D}_{val}} -\phi(v_j)^\top (V S^2 V^\top)^{-1} \phi(z_i) (1 - p^{z_i}) \\ &= \left(\sum_{v_j \in \mathcal{D}_{val}} -\phi(v_j)^\top \right) (V)^{-1^\top} S^{-2} (V)^{-1} \phi(z_i) (1 - p^{z_i}) \\ &= \left(\sum_{v_j \in \mathcal{D}_{val}} -\phi(v_j)^\top \right) V S^{-2} V^\top \phi(z_i) (1 - p^{z_i}) \\ &= \left(\sum_{v_j \in \mathcal{D}_{val}} -\nabla_\theta f(v_j, \theta)^\top \right) V S^{-2} V^\top \nabla_\theta f(z_i, \theta) (1 - p^{z_i}) \text{ (since } \mathcal{P} = I \text{ and as per Equation 8)} \\ &= \left(\sum_{v_j \in \mathcal{D}_{val}} -\nabla_\theta f(v_j, \theta)^\top \right) V S^{-2} V^\top \nabla_\theta f(z_i, \theta) (1 - p^{z_i}) \\ &= \sum_k \frac{\left(\left(\sum_{v_j \in \mathcal{D}_{val}} -\nabla_\theta f(v_j, \theta)^\top \right) v_k \right) \left(v_k^\top \nabla_\theta f(z_i, \theta) (1 - p^{z_i}) \right)}{s_{kk}^2}\end{aligned}\tag{20}$$

1006 where, v_k is the k^{th} column of V matrix and representing the k^{th} feature as per Definition 2

1007 now given the definition of the \mathcal{G}_0 and as per Equation 8, Equation 13 and Equation 15

$$\begin{aligned}\mathcal{G}_0 &= [\nabla_\theta f(z_1, \theta)^\top; \dots; \nabla_\theta f(z_n, \theta)^\top] \\ Y \mathcal{G}_0 &= U S V^\top\end{aligned}\tag{21}$$

1008 For the i^{th} training sample, this score can be further simplified by multiplying with the standard unit vector (e_i)
1009 on both sides:

$$\begin{aligned}e_i^\top Y \mathcal{G}_0 &= e_i^\top U S V^\top \\ y_i \nabla_\theta f(z_i, \theta)^\top &= u^i S V^\top\end{aligned}\tag{22}$$

1010 where u^i is a row vector associated with matrix U,

1011 multiplying both side with y_i and V we get ,

$$y_i \cdot y_i \nabla_\theta f(z_i, \theta)^\top V = y_i u^i S$$

1012 as $y_i^2 = 1$ and further multiplying both side with e_k we get

$$\begin{aligned}\nabla_\theta f(z_i, \theta)^\top V \cdot e_k &= y_i u^i S \cdot e_k \\ \nabla_\theta f(z_i, \theta)^\top v_k &= y_i u_k^i s_{kk}\end{aligned}\tag{23}$$

1013 substituting the value in Equation 20 gives :

$$|\mathcal{A}(z_i)| = \left| \sum_k \frac{\left(\sum_{v_j \in \mathcal{D}_{val}} -\nabla_\theta f(v_j, \theta)^\top \right) v_k y_i u_k^i (1 - p^{z_i})}{s_{kk}} \right|\tag{24}$$

1014 According to the given equation, for any two data points z_i and z_j where the dominant features are f_1 and f_2
1015 respectively, the contribution of these features, as per Definition 2, is represented by u_1^i and u_2^j . When both
1016 dominant features are equally represented, it follows that $u_1^i = u_2^j$ and $u_1^j < u_2^j$, $u_2^i < u_1^i$. Furthermore, if
1017 $|s_{11}| > |s_{22}|$ then as per Theorem 1 f_1 induces gradient starvation of f_2 and results in lower attribution score
1018 i.e., $|\mathcal{A}(z_i)| < |\mathcal{A}(z_j)|$. \square

G Data Annotation

G.1 Attribute Generation

We utilize ChatGPT to generate attributes for a specific dataset with the following prompt referenced from HiBug [39]. The list of attribute-value pairs generated by ChatGPT is provided in Table 9.

You are a helpful assistant to help user work on improving AI visual models. You need to discuss with your user for a description of the task that the model is working for. You need to decide if the description is complete and clear enough. The description should at least contains or infer the task object, task type, task scene. After understanding user’s task description, you should generate related visual attributes that might affect the model’s performance. You should not ask me to provide visual attributes. (Note that this is only an example visual attributes according to the previous example, do not take any of its values as default value!): “Gender , Age , Hairstyle , Hair colour” If user is satisfied with the attributes, generate the attribute form with the header formatted as “//Attribute Form//” and end with “//END//”. Attributes in the form should be splited by comma. Do not include the task object, task type, task scene. (Note that this is only an example visual attributes according to the previous example, do not take any of its values as default value!):

//Attribute Form// Gender , Age , Hairstyle , Hair colour //END//

Table 9: Details of the attribute value pair generated using ChatGPT.

Dataset	Attributes	Choices
AWA2	Size of the Animal	Small, Medium, Large, Very Large
	Fur or Skin Texture of Animals	Smooth, Rough, Furry, Scaly
	Color Pattern on Animal	Striped, Spotted, Solid Color, Mixed Colors
	Posture of Animal	Sitting, Standing, Flying, Running
	Visible Markings or Patterns	Scars, Spots, Unique Patterns
	Lighting Conditions	Bright, Dim, Natural, Artificial, Shadowy
	Background Complexity	Plain, Cluttered, Natural Habitat
	Presence of Humans	None, Nearby, Interacting
	Animal Activity State	Resting, Moving, Feeding, Playing
	Occlusions	Fully Visible, Partially Hidden
	Weather Conditions	Sunny, Cloudy, Rainy, Foggy, Snowy
	Seasonal Variations	Summer Coat, Winter Coat, Shedding Fur
CELEBA	Gender	Male, Female
	Age	Child, Teenager, Adult, Elderly
	Facial Expression	Neutral, Smiling, Frowning, Surprised
	Hairstyle	Short, Long, Bun, Braided
	Hair Color	Black, Brown, Blonde, Red
	Skin Tone	Light, Medium, Dark
	Facial Hair	Beard, Mustache, Clean-shaven
	Presence of Accessories	Glasses, Earrings, Necklace
	Lighting Conditions	Bright, Dim, Shadowed
	Makeup	Natural, Heavy, None
CIFAR-10	Size	Large, Medium, Small
	Pose/Orientation	Side View, Top View, Angled
	Lighting	Daylight, Nighttime, Shadows
	Background Complexity	Plain, Crowded
	Object Occlusion	Partially Visible, Fully Visible
GTSRB	Shape of Sign	Round, Triangular, Rectangular
	Color of Sign	Red, Blue, Yellow, White
	Size of Sign	Small, Medium, Large
	Weather Conditions	Sunny, Rainy, Foggy, Overcast
	Lighting	Daylight, Nighttime, Shadows, Glare
WaterBirds	Surrounding Environment	Forest Floor, Beach, Lake, River, Ocean, Shoreline
	Background Elements	Trees, Bushes, Rocks, Water Bodies, Sand, Human-made Structures
	Lighting Conditions	Full Daylight, Shaded Areas, Low-light, Overcast
	Weather Conditions	Sunny, Cloudy, Rainy, Foggy, Windy

G.2 Attribute-Value Annotation

We employ Llama 3.2 [78], a Vision-Language Model (VLM) with 11B parameters, to determine the most suitable value among a set of possible attributes and values for a given dataset. By iterating over a set of images in the validation set, the VLM generates metadata, which is subsequently utilized to identify the spurious features. Each image approximately takes 4-10 seconds on average to annotate, depending on the size of the image. The system prompt provided to Llama 3.2 is as follows:

You are an expert in identifying visual attributes in a given image. You will be presented with an image along with attributes and a list of choices for each of the attributes. You will be asked to choose the most suited choice for each of the attributes present in the image. Only choose one choice among all given choices

1044 *for a particular attribute. Ensure that the choice is a string. Reproduce the attribute and the choice as it is.*
 1045 *Preserve the case and the spelling. Respond with only a valid JSON object with the attributes as the keys and the*
 1046 *chosen choices as the values, and no other extra fluff. Use double inverted commas.*

1047 H Training Procedure

1048 H.1 Model Configuration and Metrics

1049 We maintained consistent hyperparameter settings across all baselines, with the only variation being the subset
 1050 of training data selected by each method. The validation set was used to identify underlying spurious biases, as
 1051 outlined in Section 3.2. For baseline comparisons, we utilized publicly available implementations. In cases where
 1052 the code was not open-sourced or experiments were not conducted on the specific datasets, we implemented
 1053 the methods and used the respective datasets for evaluation. For TracIN, we employed the fast implementation
 1054 available in the Captum library [104].

1055 Since many real-world datasets lack well-defined group structures [17], which are typically needed for evaluating
 1056 worst-group accuracy, we compare our method and baselines primarily on average accuracy. Additionally, to
 1057 understand the influence of deleting data samples in mitigating spurious features, we follow the experiment
 1058 setup defined by [28, 29, 17] and analyze the worst-case performance improvement. We used the methodology
 1059 proposed in [43] to create a subset of CELEBA with specific simplicity biases.

1060 H.2 Model Training and Datasets

1061 All experiments reported in Table 1 were conducted using the ResNet-18 architecture. The models were trained
 1062 from scratch with random initialization. For the WaterBirds dataset, the classifier was trained for 15 epochs
 1063 using stochastic gradient descent with a momentum value of 0.9 and a learning rate of 0.001. For all other
 1064 datasets, we used the Adam optimizer with a learning rate of 0.001.

1065 The AWA2-A, AWA2-B, CELEBA, models were trained for 15 epochs, while the GTSRB and CIFAR-10 models
 1066 were trained for 5 epochs. We have used the same 10 classes as mentioned in Boecking et al. [105] for all
 1067 experiments related to AWA2. For CELEBA, we used a subset of 10,000 examples from the original dataset,
 1068 with the target label being hair color (blond) and the spurious feature being gender (male). Additionally, we
 1069 induced a spurious correlation of 0.4 between the target and spurious features to mimic real-world biases. For
 1070 experiments related to ImageNet-100, we have considered the subset of the ImageNet dataset with 100 classes as
 1071 per Tian et al. [91] and trained the model for 10 epochs with the Adam optimizer. We have further considered
 1072 the attributes related to texture and shape for common classes available for the ImageNet dataset [106]. The
 1073 cutoff value to mark an attribute-value pair as spurious (τ) was decided based on the size of the corresponding
 1074 pair in the validation dataset, and the pair generating the largest difference with respect to the original dataset
 1075 was picked for analysis.

1076 To ensure a fair comparison for subset selection, we maintained uniformity in the training process across both
 1077 the original model training and the retraining process after data deletion.

1078 The experiments reported in Table 3, Table 4 were conducted using the ResNet-18 model, trained for 10 epochs
 1079 with the Adam optimizer and a learning rate of 0.001. The dataset was created by randomly sampling the
 1080 correlation factor within the range [0,1] and varying the training data size across [5000, 3000, 7000, 10000]. The
 1081 correlation attribute and target attribute were selected from the metadata provided in the CELEBA dataset [43].
 1082 Experiments on the following target-correlated attribute pairs—(arched eyebrows, receding hairline), (attractive,
 1083 mouth slightly open), (big nose, male), (goatee, bushy eyebrows), (mouth slightly open, smiling), (mouth slightly
 1084 open, wearing lipstick), (narrow eyes, eyeglasses), (pointy nose, mouth slightly open), (receding hairline, rosy
 1085 cheeks), and (male, pointy nose) are conducted with varying training dataset sizes of 3000, 5000, 5000, 5000,
 1086 5000, 5000, 7000, 7000, 7000, and 5000 samples respectively, and corresponding spurious correlation strengths
 1087 of 0.2, 0.8, 0.4, 0.4, 0.8, 0.9, 0.2, 0.6, 0.6, and 0.6 respectively. Further experiments on the target attributes
 1088 Bangs, Big Nose, Heavy Makeup, and Wearing Earrings, were conducted with correlation factors of 0.6, 0.2,
 1089 0.4, and 0.2, and with training sample sizes of 10000, 5000, 3000, and 5000, respectively. Results for these
 1090 experiments are provided in Table 14

1091 H.3 Data Attribution

1092 For the experiments reported in Table 1, approximately 3% of the data was removed from the training dataset.
 1093 We fix the data removal budget across all baselines, as it is a design choice best left to domain experts. A
 1094 smaller removal percentage prevents overpruning of the dataset (training sample for group land bird on water
 1095 is around 56 out of 4795 [29]) and highlights the precision of attribution methods by focusing on the most
 1096 harmful samples. In contrast, larger removals can obscure differences between methods due to overlapping
 1097 sample selections. For experiments related to spurious correlation in celeba, considering the stochasticity of the

training sample, we have fixed the budget size to 100 samples. Further ablation on subset size is provided in Appendix Q.1. We ensured uniformity in the data deletion process by basing it on the validation attribution score \mathcal{A} , calculated according to the respective definition of data attribution α in each baseline method, using their default hyperparameters.

For our proposed method, we performed hyperparameter tuning by selecting the rank parameter (t) from [50, 40, 10, 100] and the minimum weight (β) from [0.6, 0.7, 0.8, 0.9, 0.95]. The weight barrier (C) was chosen from [5, 10]. The optimization for Equation 10 was performed for 5000 iterations using the Adam optimizer with a learning rate of 0.0001. The value of γ is decided based on the fraction of the dataset that is removed from the training dataset. For experiments reported in Table 3 and Table 4, hyperparameter tuning was performed over the same range as in previous experiments, optimizing for both best average performance and best worst-group accuracy separately.

H.4 Textual Description

For different datasets, we used distinct textual representations of the underlying bias. The choice of textual descriptions in our experiments depends not only on the attribute-value pairs but also on the dataset itself. For instance, datasets like AWA2-A contain only label-specific information, such as color and habitat type, without an explicit attribute-value format. Therefore, a suitable textual representation for this dataset could be “*It is a (*1) animal.*” Here, (*1) represents the feature identified as a potential biased candidate. Similarly, for GTSRB, incorporating dataset context improves model performance, and a possible template could be “*(*1) of the sign is (2).*” where (*1) and (*2) are replaced by the corresponding attribute and value pair.

For datasets such as WaterBirds, AWA2-A, AWA2-B, CELEBA, GTSRB, CIFAR-10, and ImageNet-100 the textual descriptions used in the experiments related to Table 1 are provided in Table 10:

Table 10: Textual Descriptions of Spurious Feature for Different Datasets

Attribute Description	Dataset
<i>Surrounding environment in image is forest floor</i>	WaterBirds
<i>It is a domestic animal</i>	AWA2-A
<i>Size of the animal is very large</i>	AWA2-B
<i>Image of a male with blond hair</i>	CELEBA
<i>Shape of the sign is round</i>	GTSRB
<i>Size of the entity is large</i>	CIFAR-10
<i>Object has a spotted pattern</i>	ImageNet-100

For all experiments related to Table 3, Table 4, we used a standardized textual format: “*Image of a person with (*1) and (2).*” where (*1) and (*2) correspond to the target class and the correlated attribute, respectively. Further experiments using VLM-based textual description in Table 5 for the target attributes Wearing Earrings, Bangs, Big Nose, Heavy Makeup use textual description as “*Person is wearing glasses*”, “*Image of a male person*”, “*Person has long hair*”, and “*Person is wearing glasses*” respectively. For metadata, we used the same format as the Table 3.

I Comparison with Other Optimization and Data-Centric Methods

In general, ImageNet initialization [107] plays a crucial role in achieving strong worst-group accuracy. However, most of our experiments are conducted without ImageNet pretraining to better reflect practical deployment scenarios, particularly those where spurious correlations can significantly degrade model performance [107]. For a fair comparison with optimization-based methods such as gDRO [17] and JTT [16], we additionally evaluate our method on the Waterbirds dataset using a ResNet-18 model pretrained on ImageNet, along with LLM-generated attribute–value annotations. Results averaged over three independent runs are reported in Table 11. We also include comparisons with data deletion methods like D3M [11] and group-balancing approaches such as SUBG and RWG [29].

Table 11 compares the average and worst-group accuracy of our method against various robustness-based approaches on the Waterbirds dataset. Methods are grouped based on whether they require group annotations for the entire training dataset and whether they support textual bias descriptions.

Our method achieves a competitive average accuracy (0.855) and strong worst-group accuracy (0.756) without relying on group annotations, while uniquely supporting textual bias descriptions. Compared to other methods like ERM, D3M, and JTT, our method improves worst-group accuracy by +27.9% over ERM, +12% over JTT,

and +1.6% over D3M. Further comparison with D3M with the same training setup as Table 3 is provided in Table 6.

Group annotation-based methods like gDRO and RWG perform best on worst-group accuracy, but at the cost of requiring explicit group labels for the entire training dataset.

Additional challenges associated with these methods in specific applications are discussed in Section 1, Section 2.2 and Appendix E.

Table 11: Comparison of Average Accuracy and Worst group accuracy achieved by our method in comparison with other robustness-based methods on Waterbirds.

Method	Group Annotation (Train)	Supports Textual Bias Description	Average Accuracy	Worst Group Accuracy
ERM	✗	✗	0.819	0.477
D3M	✗	✗	0.903	0.740
JTT	✗	✗	0.852	0.636
Ours	✗	✓	0.855	0.756
RWG	✓	✗	0.864	0.822
SUBG	✓	✗	0.833	0.814
gDRO	✓	✗	0.886	0.836

J Empirical Validation of Theoretical Formulation

To validate our theoretical claim, we used the codebase provided by Eyuboglu et al. [43] to sample a 10k subset from CELEBA, where the attributes Male and Smiling are highly correlated. We then computed Trak scores for the training dataset using a ResNet-18 classifier trained to predict the Male label. In this setting, due to the strong correlation between Male and Smiling [43, 108], smiling may act as a spurious feature. Since the task is to distinguish males from females, we consider features like Beard and Moustache to be more causally relevant, and thus expect that samples with these features to have lower \mathcal{A} scores compared to those with Smiling.

However, statistical analysis of the detrimental attribution (\mathcal{A}) scores using T-test for the training samples reveals that Smiling has lower scores for samples compared to samples with Beard and Moustache (Table 12). The difference is statistically significant for Beard ($p < 0.001$). This supports Proposition 1, demonstrating that such effects can arise in practical scenarios.

Table 12: Mean and standard deviation of detrimental attribution ($|\mathcal{A}|$) scores for different attributes, along with statistical significance from a two-sample t-test against *Smiling*.

Attribute	Mean	Std	p-value (vs Smiling)	Significance
Smiling (spurious)	0.539	0.056	—	—
Moustache	0.545	0.044	0.1008	Not significant
Beard	0.544	0.036	0.00039	Significant

K Worst Class Performance

Table 13 reports gains in the worst-performing class for each dataset. In Awa2-A and Awa2-B, worst-class accuracy more than doubles, while in GTSRB, it improves from 50% to 70%. These results demonstrate that our method enhances class-level performance with minimal negative impact on other classes.

Table 13: Worst-class accuracy before and after retraining. The table shows the original worst-class accuracy and the corresponding value after retraining with spurious samples removed.

Dataset	Original Worst-Class Accuracy	Retrained Worst-Class Accuracy
Awa2-A	0.040	0.103
Awa2-B	0.040	0.103
CIFAR-10	0.589	0.575
GTSRB	0.500	0.700
ImageNet-100	0.100	0.100

1161 L Group-wise Accuracy Improvements

Table 14: Comparative evaluation of the proposed method (Ours) with the full training baseline (Original) and Trak, reporting the best average and best worst group accuracy (mean_{std}) across three runs.

Target Attribute	Spurious Attribute	Average Accuracy			Worst Group Accuracy		
		Original	Trak	Ours	Original	Trak	Ours
Bangs	Black Hair	0.920 _{0.007}	0.921 _{0.006}	0.923 _{0.006}	0.523 _{0.079}	0.571 _{0.053}	0.649 _{0.049}
Big Nose	Wearing Necklace	0.765 _{0.032}	0.787 _{0.009}	0.787 _{0.010}	0.127 _{0.065}	0.080 _{0.047}	0.347 _{0.148}
Heavy Makeup	Straight Hair	0.805 _{0.031}	0.800 _{0.055}	0.826 _{0.024}	0.651 _{0.137}	0.686 _{0.078}	0.716 _{0.088}
Wearing Earrings	Bags Under Eyes	0.791 _{0.020}	0.792 _{0.019}	0.798 _{0.028}	0.040 _{0.029}	0.017 _{0.029}	0.281 _{0.170}

1162 Table 15 presents group-wise accuracy before and after removing samples associated with spurious features
1163 associated with Table 14. Groups 1–4 show baseline performance, while Groups 1*–4* report results after pruning. Our method yields notable improvements in some groups without major drops in others.

Table 15: Group-wise accuracy before and after removing spurious samples. The table reports the mean accuracy and standard deviation over 3 runs. Groups 1–4 represent the training with the original dataset, while Groups 1*–4* correspond to results after data pruning.

Target Attr	Spurious-Attr	G1	G2	G3	G4	G1*	G2*	G3*	G4*
Bangs	Black Hair	0.73 _{0.07}	0.97 _{0.02}	0.53 _{0.08}	0.98 _{0.01}	0.78 _{0.06}	0.96 _{0.01}	0.59 _{0.12}	0.97 _{0.01}
Big Nose	Necklace	0.13 _{0.06}	0.94 _{0.04}	0.32 _{0.11}	0.89 _{0.07}	0.22 _{0.09}	0.94 _{0.03}	0.37 _{0.11}	0.89 _{0.06}
Heavy Makeup	Straight Hair	0.68 _{0.14}	0.80 _{0.12}	0.81 _{0.08}	0.82 _{0.07}	0.72 _{0.11}	0.83 _{0.02}	0.80 _{0.06}	0.80 _{0.02}
Earrings	Bags Under Eyes	0.09 _{0.04}	0.99 _{0.01}	0.04 _{0.03}	0.99 _{0.01}	0.28 _{0.17}	0.96 _{0.03}	0.32 _{0.16}	0.91 _{0.07}

1164

1165 M Ablation of Different Components

1166 The ablation study in Table 16 highlights the contribution of key components i.e., data Attribution and CLIP, to
1167 the overall performance of our method. For the given experiment, we have used cosine similarity with CLIP
1168 (Only CLIP) representation to remove samples that align with the description of the underlying bias. When used
1169 independently, both components provide noticeable improvements over the full training baseline, particularly
1170 in average accuracy. However, they exhibit limitations in worst-group accuracy when applied in isolation.
1171 Notably, combining both Attribution and CLIP in our full method yields the highest performance across nearly
1172 all settings, especially in worst-group accuracy, demonstrating the complementary strengths of these components
in addressing spurious correlations.

Table 16: Comparative evaluation of the proposed method (Ours) with the full training baseline (Original), Only Attribution, and Only CLIP, reporting the best average and best worst group accuracy (mean_{std}) across three runs.

Target Attribute	Spurious Attribute	Average Accuracy				Worst Group Accuracy			
		Original	Only Attribution	Only CLIP	Ours	Original	Only Attribution	Only CLIP	Ours
Bangs	Black Hair	0.920 _{0.007}	0.921 _{0.006}	0.922 _{0.008}	0.923 _{0.006}	0.523 _{0.079}	0.571 _{0.053}	0.548 _{0.074}	0.649 _{0.049}
Big Nose	Wearing Necklace	0.765 _{0.032}	0.787 _{0.009}	0.777 _{0.002}	0.787 _{0.010}	0.127 _{0.065}	0.080 _{0.047}	0.110 _{0.091}	0.347 _{0.148}
Heavy Makeup	Straight Hair	0.805 _{0.031}	0.800 _{0.055}	0.813 _{0.010}	0.826 _{0.024}	0.651 _{0.137}	0.686 _{0.078}	0.739 _{0.045}	0.716 _{0.088}
Wearing Earrings	Bags Under Eyes	0.791 _{0.020}	0.792 _{0.019}	0.791 _{0.021}	0.798 _{0.028}	0.040 _{0.029}	0.017 _{0.029}	0.009 _{0.013}	0.281 _{0.170}

1173

1174 N Architecture-based Ablation on Worst Group Accuracy and Average 1175 Accuracy

1176 We further evaluate our method on the WaterBirds dataset across different architectures, including ResNet-
1177 18, VGG16, VGG13, AlexNet, and ConvNet. Pham et al. [107] shows that the random initial weights can
1178 significantly impact the worst group performance of a model, especially in smaller networks. To replicate this
1179 setting, we tested our method under extreme conditions, maintaining consistency in textual instructions and
1180 using a single run with the same random seed across all baselines.

1181 As shown in Table 18 and Table 17, In comparison with the complete data setting our method achieves an
1182 improvement of 5.0%, 1.9%, 3.6%, 3.1%, 12.6% in worst-group accuracy for VGG16, VGG13, Convnet,
1183 ResNet18 and AlexNet architecture and an improvement of 5.2%, 7.1% 4.9% and 7.1% for VGG16, ConvNet,
1184 ResNet18, and AlexNet in average accuracy respectively.

Furthermore, compared to Trak, our method achieves an improvement of 1.1%, 2.4%, and 4.8% in average group performance for ConvNet, ResNet18, and AlexNet, respectively. Additionally, enhancements of 5.0%, 1.4%, 7.8%, 4.5%, and 12.9% in worst-group performance were observed for VGG16, VGG13, ConvNet, ResNet18, and AlexNet.

Table 17: Architecture Ablation on WaterBirds (Best Worst Group Accuracy)

Model	Original	Random	IF	TracIN	EWC	Trak	Ours
VGG16	0.053	0.050	0.064	0.064	0.062	0.053	0.103
VGG13	0.048	0.053	0.087	0.030	0.065	0.053	0.067
ConvNet	0.090	0.034	0.064	0.033	0.053	0.048	0.126
ResNet18	0.050	0.064	0.067	0.017	0.048	0.036	0.081
AlexNet	0.050	0.048	0.107	0.031	0.042	0.047	0.176

Table 18: Architecture Ablation on WaterBirds (Best Average Accuracy)

Model	Original	Random	IF	TracIN	EWC	Trak	Ours
VGG16	0.640	0.669	0.657	0.640	0.683	0.686	0.692
VGG13	0.655	0.640	0.610	0.668	0.660	0.669	0.662
ConvNet	0.654	0.705	0.640	0.721	0.711	0.714	0.725
ResNet18	0.641	0.604	0.600	0.694	0.623	0.666	0.690
AlexNet	0.644	0.650	0.586	0.693	0.658	0.667	0.715

O Experiment on Vision Transformer

Existing data attribution methods typically compute gradients over all model parameters, which often causes memory issues for large models like Vision Transformers. To address this, we follow recent works [104, 58] and calculated the gradients only for the final feature layer for both Trak and our method. However, this adaptation was incompatible with other baselines.

The results on Waterbirds for both methods are shown in Table 19.

Table 19: Best Average Accuracy and Best Worst Group performance analysis of our method in comparison with Trak and Original training of vision transformer with entire dataset.

	Average Accuracy	Worst Group Accuracy
original	0.601/0.000	0.104/0.000
Trak	0.644/0.020	0.0740/0.014
ours	0.640/0.027	0.1671/ 0.014

P Relative Comparison with the Baselines

For experiments related to Table 14, we have provided a comparison of the relative performance improvement achieved by our method against other baselines over the complete training data setting. As shown in Table 20 and Table 21, our method, on average, outperforms other baselines in terms of best average accuracy and best worst group accuracy.

Table 20: Relative improvement in Best Average Accuracy (%) achieved by our method and other baselines compared to the complete data setting(Original). The results represent the mean scores from three independent runs, with the best-performing values highlighted in **bold**.

Target Attribute	Spurious Attribute	Random	EWC	IF	TracIN	Trak	Ours
Bangs	Black Hair	-0.03	0.37	0.36	-0.07	0.16	1.05
Big Nose	Wearing Necklace	-0.07	1.66	1.10	-0.17	2.23	2.17
Heavy Makeup	Straight Hair	-0.44	-1.05	0.95	-2.83	-0.44	2.18
Wearing Earrings	Bags Under Eyes	0.3	-0.2	-1.17	-0.57	0.1	0.73

Table 21: Relative improvement in Best Worst Group Accuracy (%) achieved by our method and other baselines compared to the complete data setting(Original). The results represent the mean scores from three independent runs, with the best-performing values highlighted in **bold**.

Target Attribute	Spurious Attribute	Random	EWC	IF	TracIN	Trak	Ours
Bangs	Black Hair	6.25	15.62	6.66	7.05	4.77	12.58
Big Nose	Wearing Necklace	4.4	-3.57	5.10	4.68	-4.65	22.08
Heavy Makeup	Straight Hair	1.84	5.95	4.30	-1.24	3.55	6.54
Wearing Earrings	Bags Under Eyes	5.93	1.53	13.54	-3.46	-3.23	24.17

Q Sensitivity Analysis

Table 22 and Table 23 show the sensitivity of our proposed method on different hyperparameter values.

Table 22: Sensitivity analysis of the average accuracy of our method on the WaterBirds dataset for hyperparameters like the barrier constant (C), the matrix rank (t) (shown by rows), and the minimum weight fraction (β , shown by columns).

Barrier (C)	Rank (t)	0.6	0.7	0.75	0.8	0.85	0.9
5	40	0.657	0.619	0.648	0.673	0.650	0.640
	50	0.673	0.642	0.621	0.618	0.618	0.601
	100	0.670	0.678	0.650	0.602	0.632	0.631
10	40	0.690	0.671	0.615	0.634	0.621	0.650
	50	0.633	0.679	0.639	0.657	0.629	0.609
	100	0.653	0.663	0.602	0.642	0.660	–

Table 23: Sensitivity analysis of the worst group accuracy of our method on the WaterBirds dataset for hyperparameters like the barrier constant (C), the matrix rank (t) (shown by rows), and the minimum weight fraction (β , shown by columns).

Barrier (C)	rank (t)	0.6	0.7	0.75	0.8	0.85	0.9
5	40	0.037	0.051	0.042	0.020	0.050	0.041
	50	0.033	0.042	0.055	0.048	0.056	0.065
	100	0.020	0.036	0.041	0.081	0.041	0.050
10	40	0.009	0.037	0.056	0.051	0.044	0.030
	50	0.044	0.023	0.053	0.031	0.051	0.061
	100	0.045	0.031	0.065	0.034	0.034	–

Q.1 Performance Analysis on Different Subset Size

To further analyze model performance across different subset sizes, we conducted an ablation study where the best hyperparameters were kept fixed while varying the proportion of removed training data. The results are summarized in Table 24.

Table 24: Sensitivity analysis of the worst group accuracy and average accuracy of our method on the WaterBirds dataset for different subset sizes.

Metrics	3%	5%	15%	25%
Average Accuracy	0.69	0.645	0.682	0.712
Worst group Accuracy	0.081	0.041	0.037	0.002

R Time Taken for Subset Selection

In Figure R, we compare the time taken by our method in comparison with other baselines to select a subset of 1200 images from 60,000 images of CIFAR-10 for the instruction mentioned in Table 10. Since our method uses the attribution scores generated by Trak and improves upon it. The time taken by our method is slightly longer than Trak.

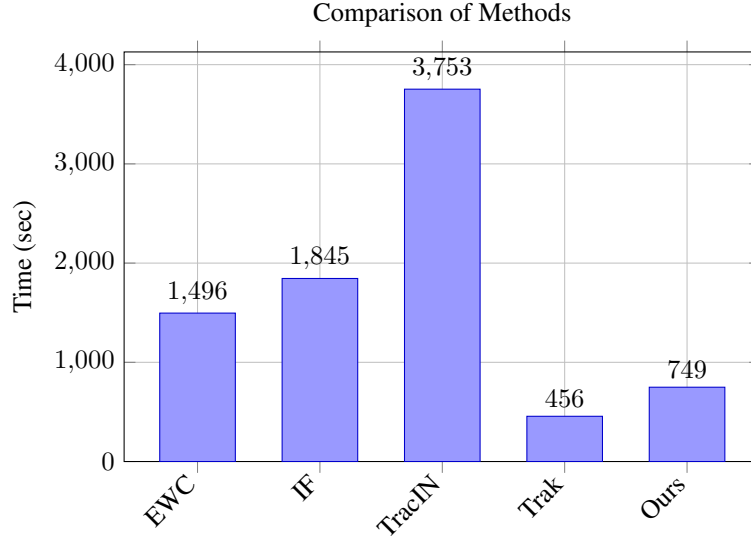


Figure 2: Comparison of time taken to select a subset of 1200 samples from a training dataset of 60,000 images of CIFAR-10 by different baselines and (Ours) for a given textual instruction.

1211 S Memory Consumption and Other Training Overhead

1212 The Table 25 reports GPU and RAM usage of our method compared to other baselines, using the same setup
 1213 described in Appendix R.

1214 As shown, our method introduces only a marginal computational overhead over Trak, which we use for computing
 1215 data attribution scores. It is to be noted that, while Trak is more memory-intensive, it produces better linear
 1216 datamodeling score (LDS) scores than other baselines [32].

Table 25: GPU and RAM utilization (in MB) of our method compared to baseline approaches.

Method	GPU Memory (MB)	RAM Usage (MB)
IF	27,749	10,578
EWC	13,221	10,520
TracIN	44,087	9,629
Trak	48,020	10,710
Ours	48,525	10,722

1217 T Workflow

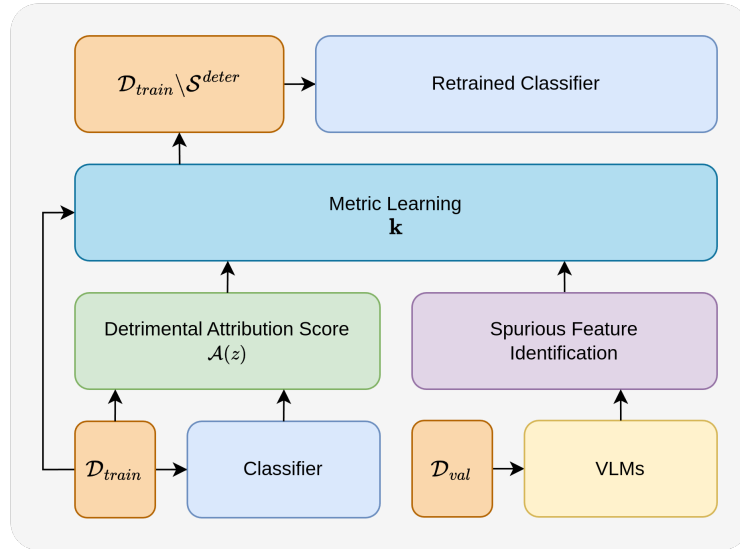


Figure 3: Diagram depicting the workflow of the proposed method

1218 U Images

1219 In this section, we have shown the images that have been removed from the training dataset. Figure 4, Figure 5,
 1220 Figure 6, and Figure 7 show the set of images that have been removed by our method from the training dataset
 1221 as (S^{deter}). For WaterBirds, GTSRB, CELEBA, and AWA2-B, respectively.

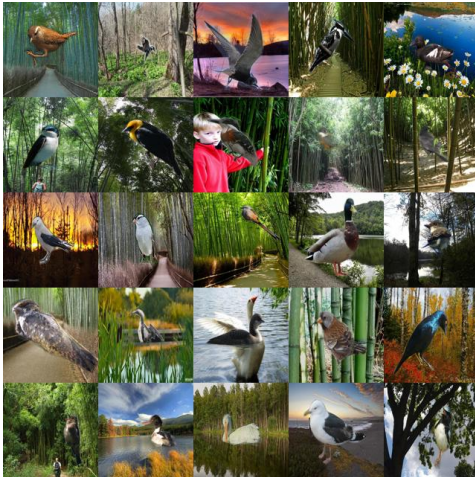


Figure 4: Set of images removed by our method for WaterBirds. The instruction set used for this experiment is “The surrounding environment in the image is forest floor”.



Figure 5: Set of Images removed by our method for GTSRB. The instruction set used for this experiment is “Shape of sign is round.”



Figure 6: Set of Images removed by our method for CELEBA. The instruction set used for this experiment is “*Image of a male with blond hair*”.

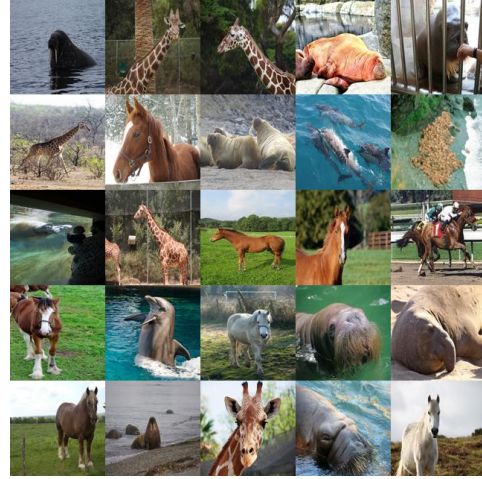


Figure 7: Set of Images removed by our method for Awa2-B. The instruction set used for this experiment is “*The size of the animal is very large.*”