LABEL-FREE MITIGATION OF SPURIOUS CORRELA-TIONS IN VLMS USING SPARSE AUTOENCODERS

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

034

035

036

037

040

041

042

043

044

046

047

051

052

Paper under double-blind review

ABSTRACT

Vision-Language Models (VLMs) have demonstrated impressive zero-shot capabilities across a wide range of tasks and domains. However, their performance is often compromised by learned spurious correlations, which can adversely affect downstream applications. Existing mitigation strategies typically depend on additional data, model retraining, labeled features or classes, domain-specific expertise, or external language models—posing scalability and generalization challenges. In contrast, we introduce a fully interpretable, zero-shot method that requires no auxiliary data or external supervision named DIAL (Disentangle, Identify, And Label-free removal). Our approach begins by filtering the representations that might be disproportionately influenced by spurious features, using distributional analysis. We then apply a sparse autoencoder to disentangle the representations and identify the feature directions associated with spurious features. To mitigate their impact, we remove the subspace spanned by these spurious directions from the affected representations. Additionally, we propose a principled technique to determine both the optimal number of spurious feature vectors and the appropriate magnitude for subspace removal. We validate our method through extensive experiments on widely used spurious correlation benchmarks. Results show that our approach consistently outperforms or matches existing baselines in terms of overall accuracy and worst-group performance, offering a scalable and interpretable solution to a persistent challenge in VLMs.

1 Introduction

Contrastive image-language models like CLIP have become foundational components in numerous applications, largely due to their remarkable zero-shot generalization capabilities Radford et al. (2021); Cherti et al. (2023). By training on web-scale data, they eliminate the need for task-specific labeled datasets, enabling efficient and scalable solutions for a wide range of downstream tasks and generative pipelines Lu et al. (2025); Zhu et al. (2025); Adila et al. (2024). However, despite strong aggregate performance, these vision-language models (VLMs) often fail on specific demographic or semantic groups, exhibiting performance far below the average Zhu et al. (2025); Chuang et al. (2023a); Yang et al. (2023). This vulnerability stems from their tendency to learn spurious correlations relying on non-causal features that are coincidentally prevalent in the training data rather than the causal task-relevant attributes Li et al. (2025). A commonly cited example in literature is where medical diagnosis predictions are being made using imaging artifacts found in the diagnostic image instead of causal disease features Lu et al. (2025); Li et al. (2025). Figure 2 shows some examples of these spurious correlations visualized through a heatmap. As these spurious correlations may not hold in real-world test data, the model's reliability and zero-shot promise are fundamentally undermined, raising serious concerns about fairness and robustness Varma et al. (2024); Chuang et al. (2023b).

In recent times, a growing body of work has sought to mitigate the spurious correlations in VLMs. Many works like Chuang et al. (2023b); Trager et al. (2023); Lauscher et al. (2020) have focused on the textual modality for debiasing, but do not address biases encoded in the visual representations. Also, methods like Lauscher et al. (2020) require domain expertise or manual specification of debiasing textual prompts. Other prominent methods Yang et al. (2023); Zhang & Ré (2022); Wang et al. (2023); Zhu et al. (2025) require fine-tuning the model or access to class and/or spurious feature labels, which negates the primary zero-shot advantage of VLMs. Recently, a few methods

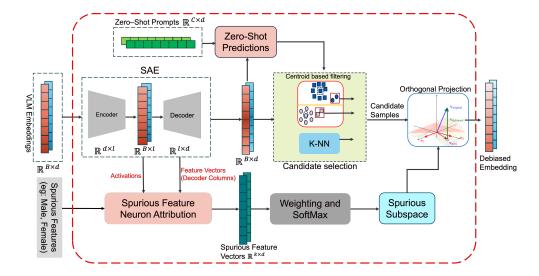


Figure 1: Overview of our proposed method. Our method takes in VLM image embeddings and spurious features of a given dataset. (e.g., "Male" and "Female" for the CelebA dataset). The entire method operates in a zero-shot setting without requiring training, external data, class labels, or spurious feature labels,.

have emerged that operate in a truly zero-shot setting Lu et al. (2025); Adila et al. (2024); Chuang et al. (2023b). However, they introduce their own set of challenges. For instance, TIE Lu et al. (2025) relies on spurious feature labels for each sample to achieve optimal performance, which are often unavailable and expensive to acquire. Moreover, although it offers a label-free variant (TIE*), both implementations practically depend on additional data to compute their scaling factors. Concurrently, methods like ROBOSHOT Adila et al. (2024) rely on Large Language Models (LLMs) to generate task-specific insights, introducing concerns about reliability, hallucination, and sensitivity to the choice of LLM Lu et al. (2025).

To address the challenges of the current methods in mitigating spurious correlations, we propose an interpretable algorithm, DIAL (Disentangle, Identify, And Label-free removal), which works in a complete zero-shot setting without requiring training, additional data, or labels (both class labels and spurious feature labels). Our framework requires two inputs: VLM embeddings of samples of a dataset and a high-level description of spurious features affecting the dataset (e.g., "Male", "Female" for CelebA). Our method unfolds in three main steps. First, guided by the insight that samples affected by spurious features often deviate from their class centroids Li et al. (2025), we identify a candidate set of potentially biased samples without class labels using zero-shot predictions as pseudo-labels. Second, we employ an off-the-shelf Sparse Autoencoder (SAE) to project these embeddings into a disentangled feature space. Within this space, we introduce a technique to reliably identify the feature directions that encode the spurious features. Finally, we debias the identified samples by removing the spurious subspace via an orthogonal projection. We also provide a technique to select the optimal parameters for our debiasing process, namely the number of spurious feature vectors (k) and the magnitude of subspace removal (λ) . The overview of our proposed approach is given in Figure 1.

We conduct extensive experiments on five standard benchmark datasets, demonstrating the efficacy of our method compared to baselines. In summary, our contributions are:

- We propose a fully zero-shot and interpretable algorithm to mitigate spurious correlations without requiring any training, additional data, class labels, or spurious feature annotations.
- We introduce a new technique to identify and isolate the spurious feature subspace from disentangled SAE representations in a zero-shot manner.
- We validate our method's effectiveness on multiple benchmarks and VLM backbones, demonstrating that our method consistently outperforms or performs comparably to current state-of-the-art baselines.

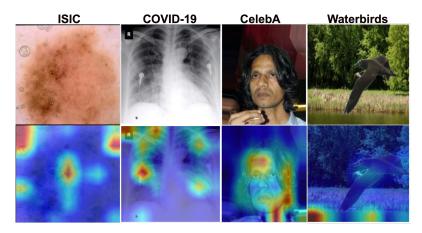


Figure 2: This figure illustrates how a CLIP model relies on spurious correlations for zero-shot predictions. For the ISIC dataset, it focuses on an image artifact instead of the lesion. For chest X-rays, it attends to a medical device rather than pneumonia indicators. On CelebA, it uses facial features instead of hair to identify 'Blond hair,' and for Waterbirds, it relies on the water background rather than the bird.

2 Related Work

Mitigation with training or labels: The problem of mitigating spurious correlations in deep learning models has been extensively studied. Techniques like Sagawa et al. (2019a); Liu et al. (2021); Yao et al. (2022); Krueger et al. (2021); Lu et al. (2024); Arjovsky et al. (2019); Idrissi et al. (2022); Yang et al. (2023); Goyal et al. (2023); Zhang & Ré (2022) aim to remove the effect of spurious correlations through reweighting the training samples, finetuning, regularization, or disparate loss functions. More recently Zhu et al. (2025) proposed to train a biased classifier to identify the group labels and debias the classifier for VLMs. Li et al. (2025) identifies the minority samples using their dispersed distribution, and learns a transformation to a bias-invariant representation. Varma et al. (2024) shows that using region-level information in the images during training helps VLMs to ignore spurious correlations. All these methods require some form of training/fine-tuning, labels, or access to the model parameters. In contrast, our method works completely in a zero-shot setting without needing any labels, fine-tuning, or access to model parameters.

Mitigation in zero-shot setting: Several of the recent works on mitigating spurious correlations in VLMs focused on doing so in a zero-shot setting. Ge et al. (2023) proposes to augment text prompts with parent and child from WordNet hierarchy to improve zero-shot generalization. Trager et al. (2023) uses the average of text prompts, which are made from combining class labels with spurious features to get debiased text prompts for each class. Dehdashtian et al. (2024) uses reproducing kernel Hilbert spaces to debias CLIP's image and text representations. Chuang et al. (2023b) proposes a closed-form method through a calibrated projection matrix to remove biased direction from clip embeddings. Lu et al. (2025) mitigates spurious correlations by translating image embeddings along the direction of spurious vectors computed from text prompts. Its main algorithm needs access to spurious feature labels for each sample, so the authors also propose a variant that adapts when spurious feature labels are not present. Additionally, both variants of TIE require access to additional data to compute the scale parameter. Additionally, both variants of TIE require access to additional data to compute the scale parameter. Additionally, both variants of the useful ones. Unlike other zero-shot approaches, our method requires no auxiliary data for parameter tuning, no spurious feature labels, and no LLM for generating insights.

Interpretable Methods for Mitigation: Some of the works have proposed using interpretability methods for mitigating spurious correlations. Wu et al. (2023) proposes an iterative framework that discovers human-interpretable spurious concepts and intervenes on training data to mitigate their influence. Chakraborty et al. (2024) uses explainability-based heatmaps for creating pseudo labels to retrain and improve robustness to spurious features in an unsupervised manner. Karvonen et al. (2024) introduces a method to evaluate an SAE based on its capacity to mitigate spurious correlations. To do this, they train linear classifiers to identify specific neurons correlated with a known spurious attribute. The activations of these identified neurons are then ablated (i.e., zeroed

out), and the resulting impact on model performance is measured. Unlike our approach, their method requires labeled training data and relies on activation zeroing rather than the removal of spurious subspace via orthogonal projection.

3 METHODOLOGY

3.1 SETUP

Let $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ be a dataset with labels $y_i \in \mathcal{Y}$. A VLM uses an image encoder ϕ_v and a text encoder ϕ_t to map inputs into a d-dimensional embedding space \mathbb{R}^d .

For zero-shot classification, a set of class prompts (e.g., "a photo of c") are tokenized and then embedded by the text encoder to produce a set of class vectors $\{p_c\}_{c=1}^{|\mathcal{Y}|}$, where $p_c = \phi_t(\text{prompt}_c)$. The probability that an image x_i belongs to class c is computed as:

$$P(y = c \mid x_i) = \operatorname{softmax}_c \left(\frac{1}{\tau} \cdot \operatorname{CosSim}(\phi_v(x_i), p_c) \right)$$

where τ is the temperature parameter.

The set of groups is defined as the $\mathcal{G} = \mathcal{Y} \times \mathcal{A}$, where \mathcal{Y} is the set of class labels and \mathcal{A} is the set of spurious attributes. We measure robustness of a VLM using three metrics: overall accuracy (Acc_{avg}) , worst-group accuracy (Acc_{wg}) , and the performance gap (Acc_{gap}) , defined as:

$$Acc_{wg} = \min_{g \in \mathcal{G}} Acc_g, \quad Acc_{gap} = Acc_{avg} - Acc_{wg}$$

The goal of our zero-shot mitigation strategy is to improve both Acc_{avg} and Acc_{wg} , and minimize Acc_{gap} , without requiring training or access to any labels.

3.2 FINDING SPURIOUS FEATURES

Our strategy is to use a pre-trained SAE to disentangle the VLM embeddings e_i and isolate feature directions corresponding to spurious attributes. An SAE decomposes an embedding into a sparse, linear combination of monosemantic features that are interpretable.

Given an embedding $e \in \mathbb{R}^d$, an SAE computes sparse feature activations $z \in \mathbb{R}^l$ and a reconstructed embedding $\hat{e} \in \mathbb{R}^d$:

$$z = act(W_{enc}e + b_{enc})$$
 $\hat{e} = W_{dec}z + b_{dec}$

Here, $W_{enc} \in \mathbb{R}^{d \times l}$ is the encoder weight matrix, and the decoder matrix $W_{dec} \in \mathbb{R}^{l \times d}$ contains the l disentangled feature vectors $\{f_j\}_{j=1}^l$ as its columns. We refer to this set of vectors as the feature dictionary, \mathcal{F} .

For each spurious attribute $a \in \mathcal{A}$ (e.g., "male" or "female"), we identify a subset of feature vectors $K_a \subset \mathcal{F}$ that strongly correlate with it. To do this, we adapt the attribution score method from Karvonen et al. (2024) to a zero-shot setting. First, we use the VLM's zero-shot classification ability to partition the reconstructed embeddings $\{\hat{e}_i\}$ from our dataset \mathcal{D} into a positive set P_a (samples exhibiting attribute a) and a negative set N_a . This is done using a prompt like "a photo of a a" and its negation.

The attribution score S for each feature vector $f_j \in \mathcal{F}$ with respect to attribute a is then calculated as:

$$S(f_j, a) = \left(\frac{1}{|P_a|} \sum_{i \in P_a} z_{i,j} - \frac{1}{|N_a|} \sum_{i \in N_a} z_{i,j}\right) \times \operatorname{CosSim}(f_j, e_a)$$

where $z_{i,j}$ is the activation of feature f_j for sample i, and $e_a = \phi_t(\text{prompt}_a)$ is the text embedding of the spurious attribute itself. This score is high when a feature's direction aligns with the attribute's semantic embedding and its activation is consistently higher for samples in the positive set.

Finally, to form the spurious feature set K_a , we select the top-k features that account for a fraction α of the total attribution mass. We sort the features f_j by $|S(f_j, a)|$ in descending order (indexed by π) and choose the smallest k such that: $\sum_{j=1}^{k} |S(f_{\pi(j)}, a)| \ge \alpha \sum_{j=1}^{l} |S(f_j, a)|$

The resulting set $K_a = \{f_{\pi(1)}, \dots, f_{\pi(k)}\}$ captures the primary directions in the embedding space associated with the spurious attribute a. The set $\mathcal{K} = \bigcup_{a \in \mathcal{A}} K_a$ contains all the feature vectors from every individual spurious feature set K_a

3.3 MITIGATING SPURIOUS FEATURES

Given the identified set of spurious feature vectors \mathcal{K} , we aim to debias the reconstructed VLM embeddings $\hat{e_i}$ by removing their components that lie in the subspace spanned by these features. To account for noise in the feature selection process, we first refine the spurious subspace by weighting each feature $f_j \in \mathcal{K}$ based on its alignment with the mean direction of the set. First, we compute the mean vector m of the spurious features: $m = \frac{1}{|\mathcal{K}|} \sum_{f_j \in \mathcal{K}} f_j$

Next, we compute a vector of alignment scores $s \in \mathbb{R}^{|\mathcal{K}|}$, where each element s_j corresponds to a feature f_j : $s_j = \beta \cdot \text{CosSim}(f_j, m)$ A weight vector w is then derived by applying the softmax function to these scores, where β is a temperature hyperparameter controlling sharpness: w = softmax(s)

To further denoise the set, we prune the features by setting weights that fall below a specified percentile to zero, yielding a filtered set of feature vectors $\mathcal{K}_f \subseteq \mathcal{K}$ with corresponding non-zero weights.

We then form a matrix V_w whose columns are the weighted feature vectors $\{w_j f_j \mid f_j \in \mathcal{K}_f\}$. We perform QR decomposition on this matrix, $V_w = QR$, to obtain an orthonormal basis Q for the refined spurious subspace. The projection of $\hat{e_i}$ onto this subspace is given by $e_{\hat{i}, \text{proj}} = QQ^T\hat{e_i}$.

The final, debiased embedding $e_{i,\hat{\text{clean}}}$ is obtained by subtracting this projection from the original embedding, scaled by a mitigation factor $\lambda \in [0,1]$: $e_{i,\hat{\text{clean}}} = \hat{e_i} - \lambda e_{i,\text{proj}}$

This procedure removes information correlated with the identified spurious concepts while preserving other essential features of the original VLM embedding.

We employ a targeted mitigation strategy, applying orthogonal projection to remove spurious features only from a subset of samples identified by our candidate selection algorithm (Alg. 1). This algorithm is designed to pinpoint samples that are likely to be affected by spurious correlations, which often lead to misclassifications.

Operating in a label-free, zero-shot setting, our approach builds on the insight from prior work Li et al. (2025) that biased samples often lie far from their true class centroid. We approximate these class centroids by using the VLM's own zero-shot predictions as pseudo-labels. To enhance the robustness of this selection against noise and outliers, we further refine the candidate set using a standard k-Nearest Neighbors (k-NN) algorithm.

Our framework has three key parameters: the number of neighbors k for k-NN, the attribution mass threshold α , and the mitigation strength λ . To select these values effectively, we propose a grid-search-based algorithm (Alg. 2) that optimizes a zero-shot score reflecting the alignment between sample embeddings and the identified spurious features.

Algorithm 1 Candidate Selection

Require: $E = \{e_i\}_{i=1}^n$, set of image embeddings. $\hat{Y} = \{\hat{y_i}\}_{i=1}^n$, set of pseudo-labels obtained from zero-shot predictions. $T = \{c \to t_c\}$, map of class labels to text embeddings. k, number of neighbors for k-NN. w, text embedding weight.

```
1: \triangleright Calculate hybrid centroids for each class c
```

```
2: for each class c \in \text{unique}(Y) do
3: \mu_c \leftarrow (1-w) \cdot \text{Mean}(\{e_i \mid \hat{y_i} = c\}) + w \cdot T[c]
```

264 4: end for

5: \triangleright Identify candidates based on centroid similarity or k-NN disagreements

```
6: M_{centroid} \leftarrow \left[\arg\max_{c'} \operatorname{CosSim}(e_i, \mu_{c'}) \neq \hat{y}_i\right]_{i=1}^n
7: M_{knn} \leftarrow \left[\text{K-NN CLASSIFY}(e_i, E, \hat{Y}, k) \neq \hat{y}_i\right]_{i=1}^n
```

 \triangleright k-NN fit with E, \hat{Y}

8: $M \leftarrow M_{centroid} \vee M_{knn}$

9: return M

4 EXPERIMENTS

4.1 DATASETS

Following the prior work Lu et al. (2025) in zero-shot spurious correlation mitigation, we use the five established benchmarks for evaluating our method. CelebA Liu et al. (2015), Waterbirds Koh et al. (2021), FMOW Christie et al. (2018) and two medical datasets ISIC Codella et al. (2019), and COVID-19 Cohen et al. (2020). All datasets except FMOW have two classes and two associated spurious features, while FMOW has 62 classes with 5 spurious features. In accordance with the prior work Lu et al. (2025); Adila et al. (2024), we define groups as a combination of class label and spurious feature. For FMOW, we define a group based on the spurious feature following the procedure given in Wu et al. (2023). For zero-shot classification, we use the same text prompts used in our prior work and evaluate all the baselines with the same text prompts. For example, for the CelebA dataset, the zero-shot text prompts we use are 'a photo of a celebrity with dark hair', and 'a photo of a celebrity with blonde hair'

4.2 Baselines

We evaluate our proposed method against existing zero-shot mitigation methods, including TIE Lu et al. (2025), ROBOSHOT Adila et al. (2024), Ideal Words Trager et al. (2023), Orth-Cali Chuang et al. (2023b), and Perception CLIP An et al. (2024). We also include the zero-shot and GroupPrompt zero-shot performance as the baselines. As established by prior works Sagawa et al. (2019b), we compare on worst group accuracy (Acc_{wg} - WG), average accuracy (Acc_{avg} - Acc), and gap between Acc and WG (Acc_{gap} - Gap). In the results, we group the baselines into two groups, one with methods that require auxiliary information through either additional data, class/spurious feature labels, or LLM insights for mitigation. This group includes Perception CLIP An et al. (2024), ROBOSHOT Adila et al. (2024), and TIE/TIE*Lu et al. (2025). The other group, which does not require any of these, is our proposed method, along with standard zero-shot, GroupPrompt classification, Ideal words Trager et al. (2023), and Orth-Cali Chuang et al. (2023b). For a fair comparison, we divide the baseline methods into these two groups in the results.

4.3 BACKBONE MODELS

Following the prior work Adila et al. (2024); Lu et al. (2025), we examine CLIP ViT-B/32 (OpenAI), and ViT-L/14 (Laion-2B) Radford et al. (2021); Cherti et al. (2023) as backbones for Waterbirds and CelebA datasets. For the FMOW dataset, we use ViT-L/14 (Laion-2B) model. For medical datasets ISIC and COVID-19 we use BiomedCLIP Zhang et al. (2023). For disentangling the representations, we use the pre-trained Matryoksha Sparse Autoencoders (MSAE) Zaigrajew et al. (2025) for all the backbone models used in the experiments. Any other SAE trained for VLMs can also be used instead of MSAE.

4.4 RESULTS

CelebA and Waterbirds:

On the CelebA dataset (results in Table 1), our method demonstrates superior performance, particularly with the ViT-B/32 backbone. It surpasses all zero-shot baselines across all three metrics, even outperforming methods that require auxiliary data, spurious feature labels, or the use of LLMs. When using the stronger ViT-L/14 backbone, our approach continues to achieve the highest worst group accuracy, lowest performance gap, underscoring its robust efficacy in mitigating spurious correlations.

For the Waterbirds dataset (results in Table 2), using the ViT-L/14 backbone, our method yields significant improvements in worst group accuracy and effectively reduces the performance gap compared to the baselines. We hypothesize that the performance on this dataset is influenced by the inherent complexity of the spurious attributes. The concepts of "land background" and "water background" are highly varied and complex, making it challenging to fully capture the corresponding feature space using only a high-level semantic description. This ambiguity may impact the preci-

Table 1: CelebA: Comparison of our mitigation method with baselines in terms of zero-shot classification. Best performance is bolded, and the second best is underlined.

Method	Setting Requirements			CLIP ViT-B/32			CLIP ViT-L/14		
	Additional Data	Class/Spurious Feature Labels	LLM	AVG (†)	WG(†)	Gap(↓)	AVG (†)	WG(†)	Gap(↓)
PerceptionCLIP	Х	Х	1	80.32	76.46	3.86	81.41	78.70	2.71
ROBÔSHOT	X	X	1	84.77	80.52	4.25	85.54	82.61	2.93
TIE	✓	✓	X	85.11	82.63	2.48	86.17	84.60	1.57
TIE*	✓	×	X	85.11	82.63	2.48	86.17	84.60	1.57
Zero-Shot	Х	Х	Х	84.27	78.89	5.38	81.20	73.35	7.85
GroupPrompt	X	X	X	80.38	74.90	5.48	77.86	68.94	8.92
Ideal words	X	X	X	80.96	78.12	2.84	89.15	76.67	12.48
Orth-Cali	X	X	X	82.31	77.92	4.39	81.39	77.69	3.70
DIAL (Ours)	X	×	X	85.54	83.47	2.17	86.87	85.24	1.63

Table 2: Waterbirds: Comparison of our mitigation method with baselines in terms of zero-shot classification. Best performance is bolded and second best is underlined.

Method	Setting Requirements			CLIP ViT-B/32			CLIP ViT-L/14		
	Additional Data	Class/Spurious Feature Labels	LLM	AVG (†)	WG(†)	Gap(↓)	AVG (↑)	WG(†)	Gap(↓)
PerceptionCLIP	Х	Х	1	82.50	59.78	22.72	86.74	54.12	32.62
ROBÔSHOT	X	X	/	71.92	54.41	17.51	64.43	45.17	19.26
TIE	✓	✓	X	79.82	71.35	8.47	84.12	78.82	5.30
TIE*	✓	×	X	76.91	61.24	15.67	78.98	61.60	17.38
Zero-Shot	Х	Х	Х	68.48	41.37	27.11	83.72	31.93	51.79
GroupPrompt	X	X	X	66.79	43.46	23.33	56.12	10.44	45.68
Ideal words	X	X	X	79.20	60.28	18.92	87.67	64.17	23.50
Orth-Cali	X	X	X	69.19	54.99	14.20	86.31	58.56	27.75
DIAL (Ours)	X	X	X	71.88	52.82	19.06	82.6	68.69	13.91

Table 3: FMOW: Comparison of our mitigation method with baselines in terms of zero-shot classification. Best performance is bolded, and the second best is underlined.

Method	Setti	ng Requirements	AVG (†)	$WG(\uparrow)$	$Gap(\downarrow)$	
1,200.00	Additional Data	Class/Spurious Feature Labels	LLM			
PerceptionCLIP	Х	Х	1	17.70	12.61	5.09
ROBÔSHOT	X	X	1	19.79	10.88	8.91
TIE	✓	✓	X	26.62	20.19	6.43
TIE*	✓	×	X	26.65	19.84	6.81
Zero-Shot	Х	Х	Х	26.02	18.06	7.96
GroupPrompt	X	X	X	14.69	8.75	5.94
Ideal words	×	X	X	20.21	11.14	9.07
Orth-Cali	X	X	X	26.11	19.45	6.66
DIAL (Ours)	×	X	X	26.09	19.90	6.19

sion of our attribution score calculation, explaining why some baselines perform better in certain configurations.

FMOW: We next evaluate our method on the challenging FMOW dataset ((results in Table 3)). Owing to the complicated nature of the dataset, following the prior work Lu et al. (2025), we use only the ViT-L/14 backbone. Our method improves over the baselines in our sub-group on the worst group accuracy while still maintaining a comparable average accuracy.

Table 4: Medical Datasets - ISIC and COVID-19: Comparison of our mitigation method with baselines in terms of zero-shot classification. Best performance is bolded, and the second best is underlined.

Method	Setting Requirements			ISIC			COVID-19		
	Additional Data	Class/Spurious Feature Labels	LLM	AVG (†)	WG(†)	Gap(↓)	AVG (†)	WG(†)	Gap(↓)
PerceptionCLIP	X	Х	/	52.74	41.55	11.19	56.87	48.84	8.03
ROBÔSHOT	X	X	/	59.84	53.30	6.54	53.10	32.75	20.35
TIE	✓	✓	X	69.90	65.87	4.03	62.50	52.17	10.33
TIE*	✓	×	X	71.68	61.11	10.57	61.08	50.22	10.86
Zero-Shot	Х	Х	Х	70.21	42.21	28.00	61.81	44.83	16.98
GroupPrompt	X	X	X	30.05	12.13	17.92	48.27	27.58	20.69
Ideal words	X	X	X	53.07	41.42	11.65	56.84	23.53	33.31
Orth-Cali	X	X	X	72.54	21.43	51.11	51.72	44.83	6.89
DIAL (Ours)	Х	×	X	70.71	68.42	2.29	<u>61.11</u>	48.28	12.83

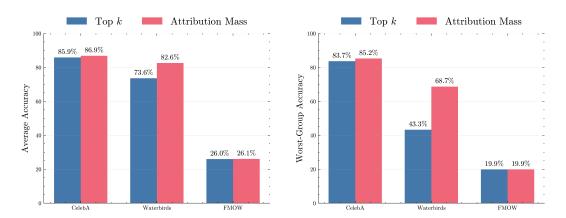


Figure 3: Comparison of spurious feature selection strategies.

Medical Datasets: The results on the medical datasets are presented in 4. On the **ISIC dataset**, our method demonstrates a substantial improvement in worst group accuracy and a corresponding reduction in the performance gap compared to all baselines. Notably, our fully zero-shot approach surpasses even those methods that rely on auxiliary data or additional labels for debiasing. Similarly, for the **COVID-19 dataset**, our approach improves over baselines in worst-group performance, it achieves this while maintaining a highly competitive average accuracy.

4.5 ABLATIONS

In this section, we justify the technical choices made in our framework through a series of empirical studies. We focus on techniques to select the optimal spurious feature vectors and removal of the spurious features. For the results reported in ablation studies, datasets CelebA, Waterbirds, and FMOW are used with Vit-L/14 as the backbone.

Selection through top k features vs attribution mass We compare the difference between selecting the top k spurious feature directions, and selecting α fraction of the attribution mass. When we run the proposed parameter search algorithm to optimize k vs alpha, we see that the latter provides better results as shown in Figure 3. This could be due to the varying representation of different features in the SAE. For example, a specific concept like "color patch" might be represented with fewer feature vectors than "land background".

Orthogonal projection vs neuron ablation Prior works have used both these techniques for concept removal. In our experiments (results shown in 4), we find that orthogonal projection is much more effective at removing the spurious features than just ablating the corresponding activations to zero. This may be attributed to orthogonal projection removing the entire spurious subspace, while ablating a specific set of neurons to zero may still leave some unidentified spurious feature vectors

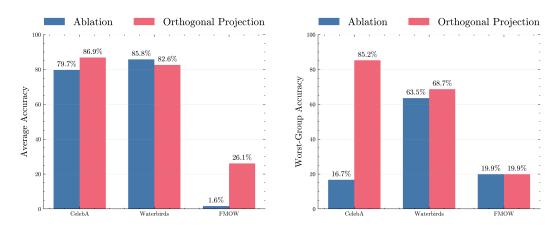


Figure 4: Comparison of spurious feature removal techniques.

watering down the mitigation. On the other hand, orthogonal projection can affect non-spurious features if they are very close to spurious features.

5 DISCUSSION

Modality and Scope: While this work focuses on mitigating spurious correlations in image embeddings, our method is modality-agnostic and can be applied to any VLM embedding. Mitigating with image modality distinguishes our approach from most zero-shot baselines that primarily target the textual modality.

SAE Parameterization: The choice of a pre-trained SAE can influence the optimal parameters (α, λ) for mitigation, as different SAEs may disentangle features at varying levels of abstraction. However, this dependency can be managed by our proposed zero-shot parameter search, which is designed to identify the optimal parameters given a set of spurious features.

Interpretable Mitigation: A key advantage of our method over prior work is its inherent transparency and interpretability. In high-stakes domains, this transparency is crucial for building trust and ensuring reliability. Our framework allows for a direct inspection of the mitigation process, providing a clear mechanism to diagnose the root causes of model failures and perform targeted debugging.

6 Conclusion

VLMs have demonstrated remarkable zero-shot capabilities, yet their performance on downstream tasks is often compromised by spurious correlations learned from web-scale training data. In this work, we introduced a novel, fully unsupervised zero-shot method to mitigate these learned biases directly within the VLM's embedding space. Our approach first employs a pre-trained SAE to get the disentangled feature representations. We then identify the directions corresponding to spurious features and remove them by applying an orthogonal projection to the VLM's image embeddings. Crucially, our method operates without requiring any additional data, training, or supervision in the form of class or spurious attribute labels. This self-contained, zero-shot nature distinguishes it from prior works that often depend on such auxiliary information or external tools like LLMs to describe spurious concepts. Furthermore, by directly manipulating image embeddings, we offer a distinct alternative to common text-embedding-based debiasing strategies. We have empirically validated our approach across five challenging datasets and multiple VLM backbones. The results demonstrate that our method consistently outperforms or performs comparably to current state-ofthe-art techniques, confirming its efficacy and potential as a practical solution for building more robust and reliable VLMs. Future work could explore multimodal debiasing, simultaneous detection, and mitigation of spurious features using SAEs.

REFERENCES

- Dyah Adila, Changho Shin, Linrong Cai, and Frederic Sala. Zero-shot robustification of zero-shot models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=fCeUoDr9Tq.
- Bang An, Sicheng Zhu, Michael-Andrei Panaitescu-Liess, Chaithanya Kumar Mummadi, and Furong Huang. More context, less distraction: zero-shot visual classification by inferring and conditioning on contextual attributes. The Twelfth International Conference on Learning Representations, 2024.
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019.
- Rwiddhi Chakraborty, Adrian Sletten, and Michael C Kampffmeyer. Exmap: Leveraging explainability heatmaps for unsupervised group robustness to spurious correlations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12017–12026, 2024.
- Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2818–2829, 2023.
- Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6172–6180, 2018.
- Ching-Yao Chuang, Varun Jampani, Yuanzhen Li, Antonio Torralba, and Stefanie Jegelka. Debiasing vision-language models via biased prompts. *arXiv preprint arXiv:2302.00070*, 2023a.
- Ching-Yao Chuang, Varun Jampani, Yuanzhen Li, Antonio Torralba, and Stefanie Jegelka. Debiasing vision-language models via biased prompts. *arXiv preprint arXiv:2302.00070*, 2023b.
- Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019.
- Joseph Paul Cohen, Paul Morrison, Lan Dao, Karsten Roth, Tim Q Duong, and Marzyeh Ghassemi. Covid-19 image data collection: Prospective predictions are the future. arXiv preprint arXiv:2006.11988, 2020.
- Sepehr Dehdashtian, Lan Wang, and Vishnu Naresh Boddeti. Fairerclip: Debiasing clip's zero-shot predictions using functions in rkhss. *arXiv* preprint arXiv:2403.15593, 2024.
- Yunhao Ge, Jie Ren, Andrew Gallagher, Yuxiao Wang, Ming-Hsuan Yang, Hartwig Adam, Laurent Itti, Balaji Lakshminarayanan, and Jiaping Zhao. Improving zero-shot generalization and robustness of multi-modal models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11093–11101, 2023.
- Sachin Goyal, Ananya Kumar, Sankalp Garg, Zico Kolter, and Aditi Raghunathan. Finetune like you pretrain: Improved finetuning of zero-shot vision models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 19338–19347, 2023.
- Badr Youbi Idrissi, Martin Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data balancing achieves competitive worst-group-accuracy. In *Conference on Causal Learning and Reasoning*, pp. 336–351. PMLR, 2022.
- Adam Karvonen, Can Rager, Samuel Marks, and Neel Nanda. Evaluating sparse autoencoders on targeted concept erasure tasks, 2024. URL https://arxiv.org/abs/2411.18895.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International conference on machine learning*, pp. 5637–5664. PMLR, 2021.

- David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In *International Conference on Machine Learning*, pp. 5815–5826. PMLR, 2021.
 - Anne Lauscher, Goran Glavaš, Simone Paolo Ponzetto, and Ivan Vulić. A general framework for implicit and explicit debiasing of distributional word vector spaces. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 8131–8138, 2020.
 - Weiwei Li, Junzhuo Liu, Yuanyuan Ren, Yuchen Zheng, Yahao Liu, and Wen Li. Let samples speak: Mitigating spurious correlation by exploiting the clusterness of samples. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 15486–15496, 2025.
 - Evan Z Liu, Behzad Haghgoo, 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*, pp. 6781–6792. PMLR, 2021.
 - 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*, pp. 3730–3738, 2015.
 - Shenyu Lu, Junyi Chai, and Xiaoqian Wang. Neural collapse inspired debiased representation learning for min-max fairness. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2048–2059, 2024.
 - Shenyu Lu, Junyi Chai, and Xiaoqian Wang. Mitigating spurious correlations in zero-shot multi-modal models. In *The Thirteenth International Conference on Learning Representations*, 2025.
 - Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.
 - 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, 2019a.
 - 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*, 2019b.
 - Matthew Trager, Pramuditha Perera, Luca Zancato, Alessandro Achille, Parminder Bhatia, and Stefano Soatto. Linear spaces of meanings: compositional structures in vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15395–15404, 2023.
 - Maya Varma, Jean-Benoit Delbrouck, Zhihong Chen, Akshay Chaudhari, and Curtis Langlotz. Ravl: Discovering and mitigating spurious correlations in fine-tuned vision-language models. *Advances in Neural Information Processing Systems*, 37:82235–82264, 2024.
 - Zhengbo Wang, Jian Liang, Ran He, Nan Xu, Zilei Wang, and Tieniu Tan. Improving zero-shot generalization for clip with synthesized prompts. *arXiv preprint arXiv:2307.07397*, 2023.
 - Shirley Wu, Mert Yuksekgonul, Linjun Zhang, and James Zou. Discover and cure: Concept-aware mitigation of spurious correlation. In *International Conference on Machine Learning*, pp. 37765–37786. PMLR, 2023.
 - Yu Yang, Besmira Nushi, Hamid Palangi, and Baharan Mirzasoleiman. Mitigating spurious correlations in multi-modal models during fine-tuning. In *International Conference on Machine Learning*, pp. 39365–39379. PMLR, 2023.
 - 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*, pp. 25407–25437. PMLR, 2022.

Vladimir Zaigrajew, Hubert Baniecki, and Przemyslaw Biecek. Interpreting clip with hierarchical sparse autoencoders. *arXiv preprint arXiv:2502.20578*, 2025.

Michael Zhang and Christopher Ré. Contrastive adapters for foundation model group robustness. *Advances in Neural Information Processing Systems*, 35:21682–21697, 2022.

Sheng Zhang, Yanbo Xu, Naoto Usuyama, Hanwen Xu, Jaspreet Bagga, Robert Tinn, Sam Preston, Rajesh Rao, Mu Wei, Naveen Valluri, et al. Biomedclip: a multimodal biomedical foundation model pretrained from fifteen million scientific image-text pairs. *arXiv preprint arXiv:2303.00915*, 2023.

Beier Zhu, Jiequan Cui, Hanwang Zhang, and Chi Zhang. Project-probe-aggregate: Efficient fine-tuning for group robustness. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 25487–25496, 2025.

A ALGORITHM

31: **return** α^* , Λ^*

594

596

598

600

601

602

603

604 605

606

607

608 609 610

611 612 613

614 615

616

617

618

619

Algorithm 2 Optimal Debiasing Parameter Search

Require: $E = \{e_i\}_{i=1}^n$, the set of original VLM embeddings. M, a boolean mask identifying the candidate subset to debias. $T_{spurious}$, a set of spurious concept text prompts (e.g., ["male", "female"]). S_{α}, S_{λ} , search ranges for hyperparameters α and λ .

Ensure: α^* , the optimal feature selection threshold. Λ^* , a map of optimal per-sample mitigation strengths for the subset.

```
1: score_{best} \leftarrow \infty
620
             2: E_{sub} \leftarrow E[M]
                                                                                    > Apply mask to get the subset of embeddings
621
             3: t_{spurious} \leftarrow \text{GetSpuriousDirection}(T_{spurious})
                                                                                                  ⊳ e.g., by averaging text embeddings
622
             4: for each \alpha \in S_{\alpha} do
623
                                                                                 \triangleright Identify the spurious subspace for the current \alpha
             5:
624
                       Q \leftarrow \text{IdentifySpuriousSubspace}(E_{sub}, \alpha)
             6:
625
             7:
                       if Q is not valid then continue
             8:
                       end if
             9:
                                \triangleright For this subspace, find the best per-sample \lambda by minimizing similarity to t_{spurious}
627
            10:
                       for each sample e_i \in E_{sub} do
628
                            \lambda_i^*, d_i^{min} \leftarrow \infty, \infty
            11:
629
                            for each \lambda \in S_{\lambda} do
            12:
630
                                 e_{i, \text{clean}} \leftarrow e_i - \lambda(QQ^Te_i)
                                                                                                                           ▶ Apply debiasing
            13:
631
                                 d_{current} \leftarrow \text{CosSim}(e_{i,\text{clean}}, t_{spurious})
                                                                                             > Score is similarity to spurious concept
            14:
632
                                 if d_{current} < d_i^{min} then
            15:
633
                                      d_i^{min} \leftarrow d_{current}\lambda_i^* \leftarrow \lambda
            16:
634
            17:
635
                                 end if
            18:
636
                            end for
            19:
637
            20:
                            \Lambda_{current}[i] \leftarrow \lambda_i^*
                            D_{min}[i] \leftarrow d_i^{min}
638
            21:
            22:
                       end for
639
            23:
                                                   \triangleright The overall score for this \alpha is the mean of the minimized similarities
640
                       score_{current} \leftarrow Mean(D_{min})
            24:
641
            25:
                       if score_{current} < score_{best} then
642
            26:
                            score_{best} \leftarrow score_{current}
643
            27:
                            \alpha^* \leftarrow \alpha
644
                            \Lambda^* \leftarrow \Lambda_{current}
            28:
645
            29:
                       end if
646
            30: end for
647
```

B REPRODUCIBILITY STATEMENT

The data pre-processing techniques we used for these experiments are the default CLIP preprocessing transforms based on the backbone architecture. All the results reported are on the test set of the datasets. MSAE models are trained with default setting mentioned in the Github repository of Zaigrajew et al. (2025) with datasets mentioned in the repo. The parameters used for our framework is can be extracted by implementing the presented parameter search algorithm.

C LLM USAGE:

We used LLMs to polish the write-up after verifying its output content. We also used LLMs precisely for searching purposes to find the relevant related works by prompting for related works based on a specific topic.