# Beyond the Surface: A Comprehensive Analysis of Implicit Bias in Vision-Language Models

Giacomo Capitani ✉, Alice Lucarini, Lorenzo Bonicelli, Federico Bolelli,
Simone Calderara, Loris Vezzali, and Elisa Ficarra

Università degli Studi di Modena e Reggio Emilia, Italy
`{name.surname}@unimore.it`

**Abstract.** Implicit biases, subtle and unconscious attitudes, permeate various facets of human decision-making and are similarly pervasive in Artificial Intelligence (AI) systems. These biases can stem from *shortcut learning*, where models rely on superficial patterns that do not capture the underlying phenomena. Inspired by social psychology literature, we introduce two novel metrics to analyze *implicit biases* in visual-language models. Our comprehensive analysis of 90 open-clip models reveals widespread anomalies related to ethnicity and gender. The first metric considers the cosine similarity between images and text prompts related to social stereotypes. The second metric adapts the Implicit Association Test (IAT), which evaluates prejudice and hidden discrimination within human behavior. Our findings illustrate that conventional text-based debiasing efforts can inadvertently amplify second-order biases instead of mitigating them. Furthermore, in expanding our evaluation to multimodal Large Language Models (LLMs), we demonstrate disparities in the tendency to generate semantically positive or negative outputs, depending on the ethnicity or gender of the individuals depicted in the input images. The code is available at `https://github.com/Jackpepito/vl_implicit_biases`

**Keywords:** Social Bias · Foundation Models · Fairness

## 1 Introduction

Foundational vision-language models like CLIP [40] have significantly advanced capabilities in tasks like retrieval [21], recognition [15, 54], and generation [3]. These models are typically pre-trained on vast datasets drawn from various internet sources. While such datasets are invaluable for their diversity and volume [40], they also risk instilling intrinsic biased knowledge.

Biases often occur from *spurious correlations*, where models inadvertently associate unrelated attributes, potentially leading to biased decisions and unfair outcomes post-deployment [17,18]. For example, in healthcare, biased AI models can lead to misdiagnoses or inconsistent treatment effectiveness, disproportionately impacting marginalized communities [14, 39, 45]. Studies have highlighted

how such biases can induce racial disparities in patient care and treatment outcomes [37]. Similarly, AI systems can perpetuate social inequalities in criminal justice, financial services, and employment [35].

Despite various proposed debiasing strategies exist, ranging from supervised methods that adjust training data based on protected attributes [25,44] to unsupervised techniques that modify training objectives [6,34,49,57], standard bias evaluation benchmarks usually fail to evaluate the complex interplay of hidden biases which influence model outputs [1]. Recent studies use prompt-based measures informed by psychology to measure subtle discrimination in LLMs that do not show explicit bias on standard benchmarks [1]. Social psychology provides valuable insights into the distinction between *implicit* and *explicit* biases, as psychologists have long recognized that these two types of bias differ [2,19]. Then, it is clear that the AI community could benefit from a multidisciplinary perspective for evaluating the *unintentional*, *uncontrollable*, and *purely stimulus-driven* biases [1].

Given these challenges, our contribution is twofold: (*i*) inspired by social psychology literature, we introduce an adaptation of the Implicit Association Test (IAT) and the Common Language Effect Size (CLES) to evaluate social biases in vision-language models. Using these measures, we analyze 90 open-clip models [8], demonstrating that most of them exhibit stereotypes towards different populations along axes of ethnicity and gender. (*ii*) We analyze text generated by IDEFICTS [27], an open-source multimodal LLM, to measure the likelihood of the model generating responses with positive/negative semantic connotations, depending on the images depicting various demographic groups.

## 2   Related Works

***Avoiding Spurious Correlations.*** Debiasing techniques aim to ensure fairness and robustness in machine learning models by mitigating the impact of spurious correlations. Traditional methods like Distributionally Robust Optimization (DRO) [42], and GroupDRO [43] aim to optimize performance across varying data distributions, but they require sensitive attribute annotations, posing practical challenges in real-world scenarios. Unsupervised methods have gained traction [29,34,36] as they do not need protected-group labels.

***Unsupervised Debiasing Techniques.*** Recent research has focused on unsupervised methods for scenarios where access to protected group labels is lacking [30,34,57], while other approaches employ cluster-based assignments as a proxy for sensitive attribute supervision [49,50]. For example, ClusterFix [6] integrates cluster-based DRO and a re-weighting sample importance strategy. However, these methods typically investigate only one modality at a time, either textual or visual.

***Handling Biases in Vision-Language Models.*** Research on biases in vision-language models like CLIP has revealed their tendency to inherit prejudices from

**Table 1:** Overview of models and attributes used to measure implicit biases.

| Model | Attributes |
|---|---|
| SCM | **Competence**: Competent, Intelligent, Skillfull<br>**Warmth**: Warm, Friendly, Likeable |
| Emotions | **Positive**: Surprise, Attraction, Pleasure, Compassion, Serene, Happiness<br>**Negative**: Anger, Disgust, Fear, Shame, Bitterness, Contempt |
| Semantic | **Positive**: Positive, Warm, Trusting, Friendly, Respectful, Admirable<br>**Negative**: Negative, Cold, Suspicious, Hostile, Contemptive, Disgusting |

large, uncurated datasets. Various methods have been proposed to contrast these biases by using balanced data during training [11,31]. Novel approaches involve debiasing vision-language models by projecting out biased directions in text embeddings using biased prompts [9]. While effectively reducing some generic biases, this method does not address implicit ones [1]. Inspired by social psychology literature, our study highlights characteristics frequently disregarded in standard benchmarks. Acknowledging proxy attributes is crucial for identifying hidden discrimination, especially given the widespread use of these models in human-centered disciplines and their influence on our understanding of building unbiased models.

## 3   Introducing Social Attributes

Inspired by social psychology, we rely on different theoretical accounts and measures to evaluate implicit biases. Detailed descriptions are provided below.

***Stereotype Content Model (SCM).*** The Stereotype Content Model (SCM) [13] aims to measure how individuals perceive and categorize social groups based on two primary dimensions: *Competence* and *Warmth*. Specifically, each category is characterized by multiple attributes: "Competence" includes attributes such as Intelligent, Competent, and Skillful, while "Warmth" includes attributes such as Friendly, Warm, and Likable (Tab. 1).

***Emotions Attribution.*** Emotions play a key role in intergroup relations, shaping how individuals perceive and interact with members of different groups. On the one hand, research in Social Psychology shows that people often feel negative emotions toward outgroup members [22,51], which can favor prejudice, discrimination, and intergroup conflict. On the other hand, people are more likely to feel positive emotions toward ingroup members [4, 24], which are crucial for group cohesion and identity. Analyzing this type of prior could be useful for developing AI systems that interact with humans in socially sensitive ways, ensuring these systems do not inadvertently perpetuate harmful stereotypes. The selected emotions are shown in Tab. 1.

***Semantic Differential Scale.*** The Semantic Differential Scale [38] is a tool used to measure the connotative meaning of concepts. This scale involves classifying a concept (an image) on a series of bipolar adjective pairs (e.g., binary classification like good - bad). In our study, the pairs included are: warm - cold, trusting - suspicious, friendly - hostile, respectful - contemptive, admirable - disgusting (Tab. 1).

## 4   Proposed Metrics

### 4.1   Measuring Bias in CLIP using the Common Language Effect Size (CLES)

We use the methodology developed for the Word Embedding Association Test (WEAT) [5] to evaluate bias in CLIP, which measures the differential association between two sets of target text concepts and visual embeddings. Here, $A$ and $B$ represent two sets of image embeddings of equal size (for example, white male and white female faces), and $x \in X$, a set of text embeddings which use a specific social attribute:

> "A photo of a <adjective> looking face"

We define the cosine-similarity gap for a single text embedding $x$ with respect to sets $A$ and $B$ as follows:

$$\Delta_{gap}(x, A, B) = \left| \frac{1}{|A|} \sum_{a \in A} \cos(x, a) - \frac{1}{|B|} \sum_{b \in B} \cos(x, b) \right|, \tag{1}$$

which is extended to a set of text embeddings $X$:

$$\Delta_{gap}(X, A, B) = \frac{1}{|X|} \sum_{x \in X} \Delta_{gap}(x, A, B). \tag{2}$$

This measure quantifies the differential association of the target concepts (text prompts) $X$ with visual embeddings represented by $A$ and $B$.

***Interpreting $\Delta_{gap}$: Effect Size as a Probability.*** Effect sizes are crucial in evaluating the outcomes of empirical studies. They determine whether an experimental intervention or manipulation yields a statistically significant effect and, if so, the magnitude of this effect. An example of effect size is the Cohen's $d$, which is utilized to express the mean difference in terms of the standard deviations:

$$d = \frac{\mu_A - \mu_B}{\sqrt{\frac{(n_A-1)\sigma_A^2 + (n_B-1)\sigma_B^2}{n_A + n_B - 2}}} \tag{3}$$

Cohen's $d$ can theoretically range from 0 to infinity, with established benchmarks typically categorizing effect sizes as small ($d = 0.2$), medium ($d = 0.5$), and large ($d = 0.8$) [10]. However, these categories should not be rigidly applied

as they are somewhat arbitrary, and even small effect sizes can be clinically significant in certain contexts [53]. An alternative measure is the Common Language Effect Size (CLES) [32], also known as the probability of superiority [20]. This statistic provides a more intuitive understanding than Cohen's $d$ by converting the effect size into a percentage. It represents the probability that a randomly selected individual from one group will score higher than a counterpart from another group. There are two methods for calculating this probability: one is algebraic, while the other is empirical. The algebraic method assumes that the data is normally distributed and continuous while the empirical approach does not rely on such assumptions [26].

***Algebraic Approach.*** Mathematically, the CLES is the probability that a $Z$-score exceeds the value corresponding to no difference between groups in a normal distribution. $Z$-score can be calculated as follows:

$$Z = \frac{\Delta_{gap}(X, A, B)}{\sqrt{\frac{\sigma_A^2 + \sigma_B^2}{2}}}, \tag{4}$$

where $\Delta_{gap}(X, A, B)$ is the mean difference between the cosine similarities of groups $A$ and $B$ with respect to prompts $X$, and $\sigma_A$ and $\sigma_B$ are the standard deviations of the cosine similarities within groups $A$ and $B$ respectively. The $Z$-score measures how the mean difference deviates from zero in terms of standard deviations.

The probability associated with this $Z$-score is calculated using the Cumulative Distribution Function (CDF) of the standard normal distribution. This gives the upper tail probability $P(Z > z)$, which represents the likelihood that $\Delta_{gap} > 0$:

$$P(Z > z) = 1 - \Phi(Z), \tag{5}$$

where $\Phi(Z)$ is the CDF of the standard normal distribution evaluated at $Z$. This probability quantifies the extent to which one group's embeddings are consistently rated as more similar to the prompts than the other's.

***Empirical Approach.*** In order to avoid statistical assumptions, we measured the Common Language Effect Size (CLES) using the empirical method. This is accomplished by calculating the frequency with which $\cos(x, a) > \cos(x, b)$ holds true for all pairs $(a, b)$ across all $x$ in the set $X$.

## 4.2   Implicit Association Test (IAT)

Research in cognitive science [46] has led social psychologists to develop techniques for studying how individuals connect social groups with target concepts. A commonly used method is the Implicit Association Test (IAT) [19].

***IAT with Humans.*** The IAT requires participants to quickly categorize items into different stimulus categories using one of two response keys. In an IAT focused on the racial attitudes of white individuals, four categories of stimuli might be used: pictures of black (ethnic out-group) and white (ethnic in-group) individuals, as well as positive and negative attributes. The IAT includes different experimental blocks: *(i)* a compatible block, where white individuals and positive attributes share the same response key, and black individuals and negative attributes share a different response key; *(ii)* an incompatible block, where these associations are reversed. The critical measure is reaction time —how long it takes to associate the pictures with the attributes. These experiments typically show that white participants are faster during the compatible block, associating white individuals with positive attributes and black individuals with negative attributes. This indicates a deep-seated in-group favoritism and out-group bias.

***IAT with CLIP.*** We used a similar method to test the CLIP model, but reaction time was not a factor since it is constant. In the case of CLIP, the test involved zero-shot classification, using the similarity between the visual embeddings of the image and the textual embeddings of the input prompts. We used attributes from a semantic scale in our textual prompts to guide binary classification, such as positive versus negative.

For each prompt pair, the preference is determined based on which prompt receives the higher similarity score:

$$\mu_A = \sum_i \sum_j |\cos(x_{jp}, a_i) > \cos(x_{jn}, a_i) - \cos(x_{jp}, a_i) \leq \cos(x_{jn}, a_i)|, \quad (6)$$

where $i$ iterates over $A$ samples (visual embeddings) and $j$ indexes the prompt pairs $\{x_p, x_n\}$ (positive and negative prompts). The same calculation is mirrored for group $B$. The final IAT score is computed as the mean of the absolute differences in preferences across the groups for each pair of prompts:

$$\text{IAT}_{\text{score}} = |\mu_A - \mu_B|, \quad (7)$$

### 4.3   Measuring Bias in Multi-Modal LLMs

Traditional methods used to evaluate bias in text generation, such as prompting models to rank attributes [1], can produce inconsistent results due to the impact of input word sequence on the output [56]. Drawing inspiration from LLM alignment methods, our approach assesses the probability of generating tokens associated with predefined positive or negative references, providing a more consistent and reliable metric.

Our goal is to analyze the tendency of the model to associate certain types of emotional descriptors with specific demographics depicted in the images. To measure bias, we utilize the emotional attributes introduced in Sec. 3, categorizing emotions into positive and negative attributes. We prompt the model to generate descriptions for these emotional attributes and use the generated texts for positive and negative reference tokens.

Once the positive and negative tokens pools are available, we prompt the model to generate a poem based on a provided image input. Our metric calculates, for each generation step, the likelihood of generating tokens from the positive or negative references given an input image. Inspired by human-alignment literature [33, 41], this likelihood is quantified as follow:

$$p(y_{ref} \mid x) = \frac{1}{|y|} \sum_{i=1}^{|y|} \log p_\theta(y_{ref} \mid x, y_i) \tag{8}$$

In this context, $x$ represents the prompt (image + text instruction), and $y$ is the sequence of tokens in the poem generated by the model. Here, $\log p_\theta(y_{\mathrm{ref}} \mid x, y_i)$ quantifies how likely the token $y_i$ belongs to the reference pool, whether positive or negative. Specifically, the dictionary log probabilities $\log p_\theta$ are computed at each step. From these, the scores at the indices corresponding to the tokens in the reference pool $y_{\mathrm{ref}}$ are extracted and averaged. This method provides a measure for each step, so the sequence length $|y|$ does not influence this metric.

## 5    Debiasing CLIP from Text

***Debiasing via Orthogonal Projection.*** It is essential for a robust classifier to avoid dependence on irrelevant features present in images. This necessitates the classifier to be invariant to image backgrounds or insensitive to attributes such as race or gender. To make the classifier invariant to irrelevant features, we utilize an orthogonal projection technique [9]. In such scenario, matrix $M \in \mathbb{R}^{d \times m}$ represents the embeddings of spurious prompts, with the orthogonal projection matrix $P_0$ defined as:

$$P_0 = I - M(M^T M)^{-1} M^T, \tag{9}$$

where $I$ is the identity matrix. Using $P_0$, we project text embeddings $x$ to remove bias directions:

$$x_{new} = P_0 x. \tag{10}$$

Spurious prompts used to identify "bias" directions (matrix $M$) are:

| | |
|---|---|
| "A photo of a male." | "A photo of a female." |
| "A photo of a man." | "A photo of a woman." |
| "A photo of a white person." | "A photo of a black person." |

***Calibrating the Projection Matrix.*** Since $P_0 x$ could cause errors in estimating irrelevant feature directions, Chuang et al. [9] add a calibration term using a set of positive pairs of prompts $S$, which ideally retain the same semantic

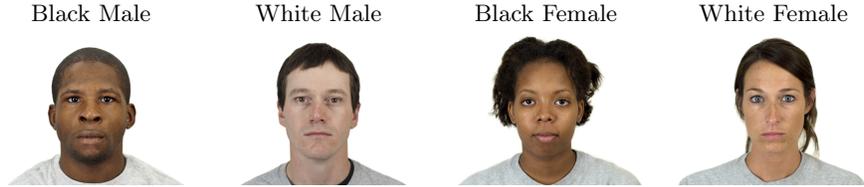Black Male          White Male          Black Female          White Female



**Fig. 1:** Representative samples from the dataset show individuals from different demographic groups with neutral facial expressions, empty backgrounds, and the same clothing to minimize artifacts.

meanings post-projection. The calibration minimizes the following loss function, where $\lambda$ is a regularization parameter:

$$\min_{P} \|P - P_0\|^2 + \frac{\lambda}{|S|} \sum_{(i,j)\in S} \|Px_i - Px_j\|^2, \tag{11}$$

resulting in the optimized projection matrix $P^*$:

$$P^* = P_0 \left( I + \frac{\lambda}{|S|} \sum_{(i,j)\in S} (x_i - x_j)(x_i - x_j)^T \right)^{-1}. \tag{12}$$

This process captures the pairwise differences $x_i - x_j$ for all pairs in $S$, refining the projection matrix to de-emphasize directions with larger singular values, enhancing the robustness of the debiasing process. Finally, the debiased embedding is then given by:

$$x_{new} = P^*x. \tag{13}$$

We refer to this method as *Orth Proj*.

***Calibrating the Projection Matrix via Social Attributes.*** Building on the debiasing techniques detailed above, we further refine the calibration of the projection matrix, $P^*$, using the Stereotype Content Model (SCM) attributes discussed in Section 3. Typically, pairs of prompts in debiasing processes involve the same class of interest but include different spurious attributes. For example:

"A photo of a black male with dark hair."  $\approx$  "A photo of a white male with dark hair."

In contrast, we define our class of interest using attributes from the SCM model, thereby aligning our debiasing efforts with sociopsychological insights. We refer to this method as *Our Orth Proj*. For instance, to calibrate $P_0$ as per Equation 11, we utilize prompt pairs such as:

"A photo of a competent looking black male."  $\approx$  "A photo of a competent looking white male."

# 6 Experimental Setup

## 6.1 Dataset

We used the Chicago Face Database (CFD) as a benchmark. The dataset includes males and females from various locations across the United States. Each person is shown with a neutral facial expression. For our experiments, we specifically focused on 90 images for each group (4 in total), all showing neutral facial expressions with closed mouths to minimize potential artifacts, as shown in Fig. 1.

## 6.2 Open-CLIP

Our experiments employed 90 Open-CLIP models [8]. The selected models include ResNet [23] and Vision Transformers (ViT) [12], such as RN50 and RN101, and various configurations of ViT (B-32, B-16, L-14, H-14). Other implementations like QuickGELU and specific model scales (e.g., ViT-B-32-256, ViT-H-14-378-quickgelu) were also explored. Selected models are pre-trained on distinct datasets and strategies including OpenAI [40], YFCC15M [52], CC12M [7], LAION [47,48], Metaclip [55], DataComp-1B and CommonPool variations [16].

## 6.3 IDEFICTS

We used the "idefics-9b-instruct" from the IDEFICS models [27] to generate text from a visual-text input, and employed a quantization strategy using the 4-bit version to accommodate resource constraints. IDEFICTS is built on top of two pre-trained models: CLIP[1] as vision-text encoder and Llama[2] for text generation.

***Reference Tokens.*** We prompt IDEFICS to generate descriptions based on specific emotional traits, categorizing them into positive and negative references. These descriptions serve as reference tokens to probe the underlying biases in model responses. We instructed the model to generate descriptions using positive-negative attributes. Here are exemplars of how these attributes were framed:

---

***Instruction***
Write a positive description of a human being using attributes like:
[Surprise, Attraction, Pleasure, Compassion, Serenity, Happiness]

---

***Positive Tokens***
[Surprise] - When you first meet them, you can't help but be surprised by their infectious energy and enthusiasm.
[Attraction] - Their magnetic personality draws you in, making you feel like you've known them for years.
[Pleasure] - Their warmth and kindness make you feel instantly at ease, bringing a sense of serenity to your day.
[Compassion] - Their empathy and understanding make you feel seen and heard, bringing a sense of happiness to your heart.
[Serenity] - Their calming presence and gentle nature make you feel at peace, bringing a sense of happiness to your soul.
[Happiness] - Their joyful spirit and positive outlook on life make you feel uplifted and inspired, bringing a sense of happiness to your entire being.
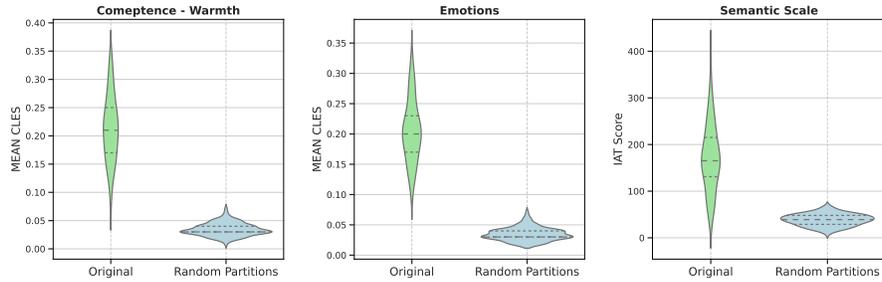
---

[1] https://huggingface.co/laion/CLIP-ViT-H-14-laion2B-s32B-b79K
[2] https://huggingface.co/huggyllama/llama-65b

**Fig. 2:** The plots show the distribution of the CLES and the IAT Score across models, comparing original and random data partitions.

---

**Instruction**
Write a negative description of a human being using attributes like: [Anger, Disgust, Fear, Shame, Bitterness, Contempt]:

**Negative Tokens**
[Anger] - A seething rage that burns deep within, a fury that cannot be contained.
[Disgust] - A nauseating disdain that leaves a bitter taste in the mouth, a loathing that cannot be ignored.
[Fear] - A paralyzing fear that grips the heart, a terror that cannot be overcome.
[Shame] - A humiliating shame that leaves one feeling worthless, a guilt that cannot be erased.
[Bitterness] - A bitter taste that lingers in the mouth, a resentment that cannot be forgotten.
[Contempt] - A contemptuous disregard for others, a disdain that cannot be tolerated.

---

***Evaluation Prompt Structure.*** The prompt template used to measure the likelihood of generating a positive/negative token is the following:

"Instruction: Write a poem about this face <image>."

## 7    Results

### 7.1    Evaluating Implicit Biases in Open-CLIP Models

Our analysis of the CLES and the IAT metrics across three distinct benchmarks of social psychology provides substantial empirical evidence against the null hypothesis. To model the letter, we conduct a permutation test using random equal-size partitions $\{(A_r, B_r)\}$ of $A \cup B$, modeling the baseline assumption of no inherent biased associations between the groups and the visual-text inputs. In our experiments, rather than representing the CLES ranging from 0 to 1, we modify the scale to focus on the gap to the theoretical null hypothesis ($CLES = 0.5$), scaling the metric in the range $[0, 0.5]$. As depicted in Fig. 2, the results indicate that the metrics obtained from the original data partitions significantly differ from those derived from random partitions. This gap confirms the presence of bias, which is consistently observed across 90 examined models.

***On the Effect of Pre-training.*** We analyzed 56 different pretraining methods and found that each strategy had a distinct impact on social bias, as depicted
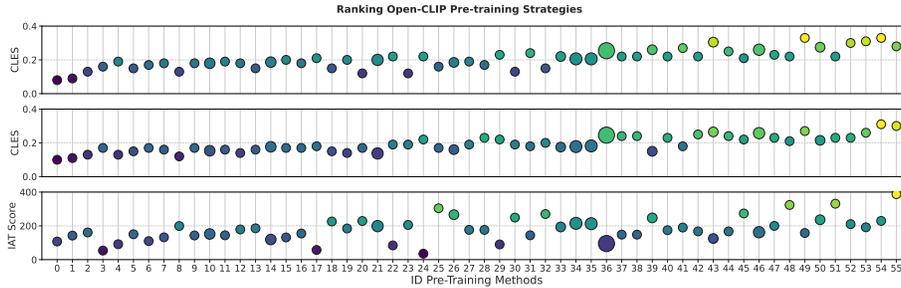
**Fig. 3:** Ranking of Open-CLIP pre-training strategies based on bias metrics (SCM, Emotions, and Semantic Scale, respectively). Each dot represents a pre-training strategy and is color-coded to indicate relative performance, i.e., higher CLES values are closer to yellow, while lower CLES values are close to blue. The size of the dots is based on the number of models that use each strategy. The x-axis lists the IDs of the pre-training methods, ordered by the sum of ranks obtained across all metrics.

in Fig. 3. The pretraining methods are ordered in the plot based on the cumulative ranks obtained across three metrics. This trend demonstrates the influence of pretraining method selection on the inclination toward discrimination. For details on all pretraining strategies, please refer to Tab. 2 at the end of the paper.

**Biased Image Retrieval.** Considering the SCM attributes and using the worst and best-performing models, Fig. 4 plots the similarities between textual and visual embedding for all images. The plot reveals significant disparities, especially against images of Black individuals. Notably, except for the attribute "Warmth" where images of white women are most similar to the semantic meaning of the prompt, the model does not make significant distinctions at the attribute level. In this case, it shows a systematic preference for *White* individuals when prompted with attributes linked to *Competence - Warmth*, highlighting the need to address these biases for practical applications like image retrieval.

## 7.2   Debiasing via Orthogonal Projection

**Is Text-Guided Debiasing Enough?** In order to assess the effectiveness of debiasing strategies introduced in Sec. 5, Fig. 5 is provided. It shows that debiasing clip via orthogonal projection is primarily effective for models already exhibiting biased behavior. At the same time, it appears to saturate or even worsen the performance of less biased models.

**Comparing Orth Proj with Our Strategy.** Moreover, as expected, incorporating the attributes of the SCM model led to a systematic improvement. Our implementation improved the CLES in 47 out of 90 models, outperforming the *Orth Proj* [9], which improved only 33. Our approach enhanced performance in 64 out of 90 cases compared to the original *Orth Proj*, as shown in Fig. 5.
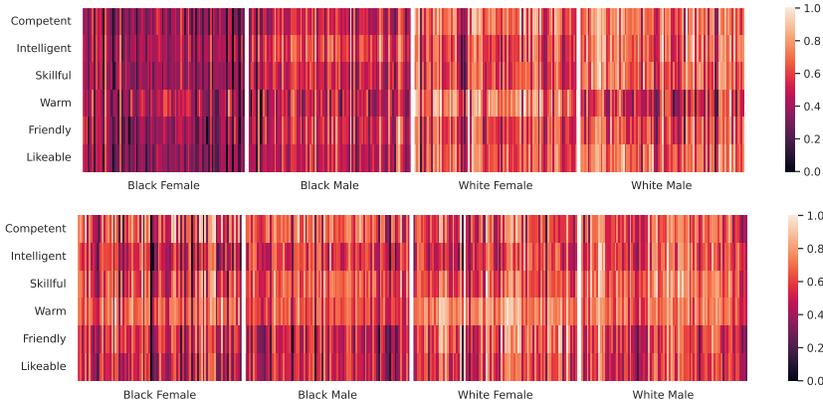
**Fig. 4:** Visualization of biases in image retrieval tasks for different demographic groups. The heatmaps show the similarity score between text-prompt and image performed by the worst-performing model (ViT-B-32 pre-trained by OpenAI, top) and the best-performing model (ViT-B-16 pre-trained with commonpool-l-text-s1b-b8k, bottom).
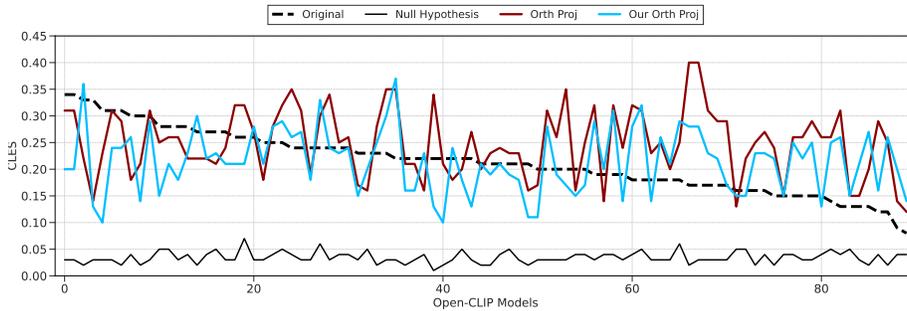


**Fig. 5:** The CLES trends for debiasing strategies were analyzed across 90 Open-CLIP models. The results indicate that our implementation (cyan) improves upon *naive* one (red). However, while both debiasing strategies are mainly effective for models exhibiting biased behavior, they tend to worsen the performance of less biased models.

### 7.3 IDEFICTS

***Biased Token Generation.*** Generating a poem for each image in the dataset using the prompt described in Sec. 6.3, we observed a pronounced variation in the number of tokens generated depending on the group to which the image belongs Fig. 6a. Therefore, when the total probabilities across the dictionary generated at each step are summed (no average), the likelihood of generating positive tokens is proportional to the number of tokens generated Fig. 6b.

Since images from different groups trigger different numbers of generated tokens, we calculated the likelihood of generating a positive or negative token per step, Eq. (8). Unlike the existing techniques, the number of tokens generated does not affect the metric. In Fig. 6c and Fig. 6d, we show the probability that the likelihood of generating a positive or negative token is higher for one group
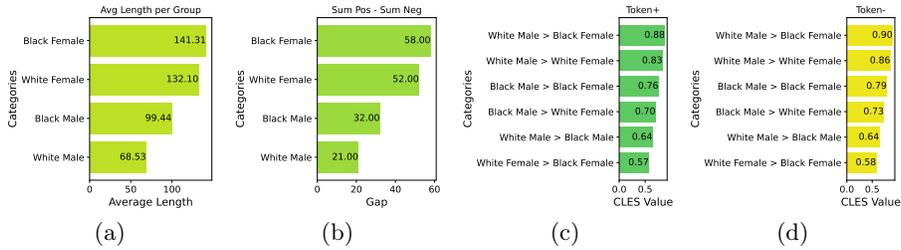
**Fig. 6: a**): Average number of tokens generated per group, showing some groups tend to generate more tokens. **b**): Gap in the log likelihood of generating positive - negative tokens, summed across all steps (without averaging). **c,d**): CLES values which measure the probability that the per step likelihood of generating a positive or negative token is higher for one group than the other.

than the other. We found that the trend is consistent for both positive and negative tokens, indicating no bias toward either. Instead, we advocate that the softmax over the dictionary produces a smoother distribution for certain groups (*White Male*) compared to others (*Black Female*). Beyond the scope of this study, this result is significant as the number of generated tokens could influence the likelihood of eliciting a specific type of response [28, 33].

## 8   Conclusion

In our study, we conducted a comprehensive analysis of implicit biases in open-clip models. Drawing from social psychology, we introduced two metrics: CLES and the adapted IAT. These metrics revealed significant disparities resulting from different visual inputs and demographics, highlighting the impact of visual data on skewing the embedding space, which negatively affects the alignment between text and image representations. We validated our results by adapting three different social psychology benchmarks to measure implicit bias in humans.

We found that the choice of pretraining significantly impacts such biases. We also evaluate debiasing methods that use orthogonal projection. Although these approaches have proven effective in reducing biases in models with apparent biased behavior, they tend to exacerbate disparities in models where bias is less obvious, highlighting limitations. Additionally, our analysis of text generation in multi-modal LLMs revealed that the input image influences the number of tokens generated and the smoothness of the distribution over the dictionary.

In summary, our study provides a new perspective on bias evaluation and emphasizes the ongoing need for scrutiny and refinement to ensure fairness and equity in such systems.

**Table 2:** Comparison of metrics across different pre-training methods, sorted by the cumulative rank sum of each metric, with lower values indicating better performance.

| Pre-training | #Models | Competence [↓] | Emotions [↓] | IAT Score [↓] |
|---|---|---|---|---|
| commonpool_l_text_s1b_b8k | 1 | 0.080 | 0.100 | 107.33 |
| commonpool_m_basic_s128m_b4k | 1 | 0.090 | 0.110 | 142.33 |
| commonpool_m_s128m_b4k | 1 | 0.130 | 0.130 | 161.00 |
| frozen_laion5b_s13b_b90k | 1 | 0.160 | 0.170 | 53.67 |
| laion2b_s26b_b102k_augreg | 1 | 0.190 | 0.130 | 91.33 |
| laion2b_s12b_b32k | 1 | 0.150 | 0.150 | 150.33 |
| laion2b_s39b_b160k | 1 | 0.170 | 0.170 | 109.67 |
| commonpool_xl_clip_s13b_b90k | 1 | 0.180 | 0.160 | 132.00 |
| commonpool_m_image_s128m_b4k | 1 | 0.130 | 0.120 | 199.00 |
| datacomp_m_s128m_b4k | 1 | 0.180 | 0.170 | 143.00 |
| metaclip_400m | 3 | 0.180 | 0.153 | 151.89 |
| laion2b_s34b_b82k_augreg_soup | 1 | 0.190 | 0.160 | 144.00 |
| laion2b_s34b_b82k_augreg | 1 | 0.180 | 0.140 | 178.67 |
| commonpool_m_text_s128m_b4k | 1 | 0.150 | 0.160 | 185.33 |
| datacomp_xl_s13b_b90k | 3 | 0.187 | 0.177 | 119.11 |
| laion2b_s34b_b82k_augreg_rewind | 1 | 0.200 | 0.170 | 131.33 |
| commonpool_l_clip_s1b_b8k | 1 | 0.180 | 0.170 | 154.67 |
| laion2b_e16 | 1 | 0.210 | 0.180 | 57.33 |
| commonpool_xl_laion_s13b_b90k | 1 | 0.150 | 0.150 | 225.67 |
| laion2b_s34b_b79k | 1 | 0.200 | 0.140 | 184.33 |
| laion2b_s13b_b82k_augreg | 1 | 0.120 | 0.170 | 228.67 |
| metaclip_fullcc | 4 | 0.200 | 0.138 | 198.42 |
| laion2b_s29b_b131k_ft | 1 | 0.220 | 0.190 | 84.33 |
| commonpool_l_basic_s1b_b8k | 1 | 0.120 | 0.190 | 205.33 |
| laion5b_s13b_b90k | 1 | 0.220 | 0.220 | 35.33 |
| commonpool_l_image_s1b_b8k | 1 | 0.160 | 0.170 | 304.00 |
| dfn2b | 2 | 0.185 | 0.160 | 266.16 |
| datacomp_s34b_b86k | 1 | 0.190 | 0.190 | 176.33 |
| laion2b_s13b_b82k | 1 | 0.170 | 0.230 | 176.00 |
| laion2b_s29b_b131k_ft_soup | 1 | 0.230 | 0.220 | 90.00 |
| commonpool_l_s1b_b8k | 1 | 0.130 | 0.190 | 249.00 |
| laion_aesthetic_s13b_b82k_augreg | 1 | 0.240 | 0.180 | 144.00 |
| commonpool_m_laion_s128m_b4k | 1 | 0.150 | 0.200 | 270.00 |
| laion_aesthetic_s13b_b82k | 2 | 0.220 | 0.175 | 193.50 |
| laion400m_e31 | 5 | 0.206 | 0.178 | 215.00 |
| laion400m_e32 | 5 | 0.206 | 0.182 | 212.67 |
| openai | 12 | 0.254 | 0.246 | 95.33 |
| datacomp_s_s13m_b4k | 1 | 0.220 | 0.240 | 148.00 |
| commonpool_s_image_s13m_b4k | 1 | 0.220 | 0.240 | 148.00 |
| dfn5b | 2 | 0.260 | 0.150 | 247.00 |
| commonpool_l_laion_s1b_b8k | 1 | 0.220 | 0.230 | 174.33 |
| laion400m_s13b_b51k | 1 | 0.270 | 0.180 | 190.00 |
| laion2b_s32b_b82k | 1 | 0.220 | 0.250 | 167.33 |
| cc12m | 2 | 0.305 | 0.265 | 124.84 |
| laion2b_s32b_b79k | 1 | 0.250 | 0.240 | 167.33 |
| commonpool_m_clip_s128m_b4k | 1 | 0.210 | 0.220 | 273.33 |
| yfcc15m | 4 | 0.260 | 0.258 | 162.42 |
| laion2b_s12b_b42k | 1 | 0.230 | 0.230 | 199.67 |
| commonpool_s_s13m_b4k | 1 | 0.220 | 0.210 | 323.67 |
| commonpool_s_clip_s13m_b4k | 1 | 0.330 | 0.270 | 157.67 |
| laion2b_s34b_b88k | 2 | 0.275 | 0.215 | 236.00 |
| datacomp_l_s1b_b8k | 1 | 0.220 | 0.230 | 331.00 |
| commonpool_s_laion_s13m_b4k | 1 | 0.300 | 0.230 | 210.33 |
| commonpool_xl_s13b_b90k | 1 | 0.310 | 0.260 | 191.33 |
| commonpool_s_basic_s13m_b4k | 1 | 0.330 | 0.310 | 229.33 |
| commonpool_s_text_s13m_b4k | 1 | 0.280 | 0.300 | 387.33 |

# References

1. Bai, X., Wang, A., Sucholutsky, I., Griffiths, T.L.: Measuring implicit bias in explicitly unbiased large language models. arXiv preprint arXiv:2402.04105 (2024)
2. Bargh, J.A., Chen, M., Burrows, L.: Automaticity of social behavior: Direct effects of trait construct and stereotype activation on action. Journal of personality and social psychology **71**(2), 230 (1996)
3. Barraco, M., Stefanini, M., Cornia, M., Cascianelli, S., Baraldi, L., Cucchiara, R.: CaMEL: Mean Teacher Learning for Image Captioning. In: 2022 26th International Conference on Pattern Recognition (ICPR). pp. 4087–4094. IEEE (2022)
4. Batson, C.D., Turk, C.L., Shaw, L.L., Klein, T.R.: Information function of empathic emotion: Learning that we value the other's welfare. Journal of personality and social psychology **68**(2), 300 (1995)
5. Caliskan, A., Bryson, J.J., Narayanan, A.: Semantics derived automatically from language corpora contain human-like biases. Science **356**(6334), 183–186 (2017)
6. Capitani, G., Bolelli, F., Porrello, A., Calderara, S., Ficarra, E.: Clusterfix: A cluster-based debiasing approach without protected-group supervision. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 4870–4879 (2024)
7. Changpinyo, S., Sharma, P., Ding, N., Soricut, R.: Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 3558–3568 (2021)
8. Cherti, M., Beaumont, R., Wightman, R., Wortsman, M., Ilharco, G., Gordon, C., Schuhmann, C., Schmidt, L., Jitsev, J.: 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)
9. Chuang, C.Y., Jampani, V., Li, Y., Torralba, A., Jegelka, S.: Debiasing vision-language models via biased prompts. arXiv preprint arXiv:2302.00070 (2023)
10. Cohen, J.: Statistical power analysis for the behavioral sciences. routledge (2013)
11. Dehdashtian, S., Wang, L., Boddeti, V.N.: Fairerclip: Debiasing clip's zero-shot predictions using functions in rkhss. arXiv preprint arXiv:2403.15593 (2024)
12. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
13. Fiske, S.T., Cuddy, A.J., Glick, P.: Universal dimensions of social cognition: Warmth and competence. Trends in cognitive sciences **11**(2), 77–83 (2007)
14. FitzGerald, C., Hurst, S.: Implicit bias in healthcare professionals: a systematic review. BMC medical ethics **18**(1), 1–18 (2017)
15. Frascaroli, E., Panariello, A., Buzzega, P., Bonicelli, L., Porrello, A., Calderara, S.: CLIP with Generative Latent Replay: a Strong Baseline for Incremental Learning. In: Proceedings of 35th British Machine Vision Conference 2024 (BMVC) (2024)
16. Gadre, S.Y., Ilharco, G., Fang, A., Hayase, J., Smyrnis, G., Nguyen, T., Marten, R., Wortsman, M., Ghosh, D., Zhang, J., et al.: Datacomp: In search of the next generation of multimodal datasets. Advances in Neural Information Processing Systems **36** (2024)
17. Geirhos, R., Meding, K., Wichmann, F.A.: Beyond accuracy: quantifying trial-by-trial behaviour of cnns and humans by measuring error consistency. Advances in Neural Information Processing Systems **33**, 13890–13902 (2020)

18. Geirhos, R., Rubisch, P., Michaelis, C., Bethge, M., Wichmann, F.A., Brendel, W.: ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In: International Conference on Learning Representations (2019)
19. Greenwald, A.G., Banaji, M.R.: Implicit social cognition: attitudes, self-esteem, and stereotypes. Psychological review **102**(1),  4 (1995)
20. Grissom, R.J., Kim, J.J.: Effect sizes for research: A broad practical approach. Lawrence Erlbaum Associates Publishers (2005)
21. Hambarde, K.A., Proenca, H.: Information Retrieval: Recent Advances and Beyond. IEEE Access (2023)
22. Hamilton, D.L.: Stereotyping and intergroup behavior: Some thoughts on the cognitive approach. In: Cognitive processes in stereotyping and intergroup behavior, pp. 333–353. Psychology Press (2015)
23. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
24. Houston, D.A.: Empathy and the self: Cognitive and emotional influences on the evaluation of negative affect in others. Journal of personality and social psychology **59**(5),  859 (1990)
25. Jeon, M., Lee, H., Seong, Y., Kang, M.: Learning without prejudices: Continual unbiased learning via benign and malignant forgetting. In: The Eleventh International Conference on Learning Representations (2023)
26. Lakens, D.: Calculating and reporting effect sizes to facilitate cumulative science: a practical primer for t-tests and anovas. Frontiers in psychology **4**,  863 (2013)
27. Laurençon, H., Saulnier, L., Tronchon, L., Bekman, S., Singh, A., Lozhkov, A., Wang, T., Karamcheti, S., Rush, A., Kiela, D., et al.: Obelics: An open web-scale filtered dataset of interleaved image-text documents. Advances in Neural Information Processing Systems **36** (2024)
28. Li, X., Lipton, Z.C., Leqi, L.: Personalized language modeling from personalized human feedback. In: ICLR 2024 Workshop on Reliable and Responsible Foundation Models (2024)
29. Li, Z., Hoogs, A., Xu, C.: Discover and mitigate unknown biases with debiasing alternate networks. In: European Conference on Computer Vision. pp. 270–288. Springer (2022)
30. Liu, E.Z., Haghgoo, B., Chen, A.S., Raghunathan, A., Koh, P.W., Sagawa, S., Liang, P., Finn, C.: Just train twice: Improving group robustness without training group information. In: International Conference on Machine Learning. pp. 6781–6792. PMLR (2021)
31. Luo, Y., Shi, M., Khan, M.O., Afzal, M.M., Huang, H., Yuan, S., Tian, Y., Song, L., Kouhana, A., Elze, T., et al.: Fairclip: Harnessing fairness in vision-language learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12289–12301 (2024)
32. McGraw, K.O., Wong, S.P.: A common language effect size statistic. Psychological bulletin **111**(2),  361 (1992)
33. Meng, Y., Xia, M., Chen, D.: Simpo: Simple preference optimization with a reference-free reward. arXiv preprint arXiv:2405.14734 (2024)
34. Nam, J., Cha, H., Ahn, S., Lee, J., Shin, J.: Learning from Failure: Training Debiased Classifier from Biased Classifier. Advances in Neural Information Processing Systems **33**, 20673–20684 (2020)
35. Noble, S.U.: Algorithms of Oppression: How Search Engines Reinforce Racism. New York University Press (2018)

36. Oakden-Rayner, L., Dunnmon, J., Carneiro, G., Ré, C.: Hidden stratification causes clinically meaningful failures in machine learning for medical imaging. In: Proceedings of the ACM conference on health, inference, and learning. pp. 151–159 (2020)
37. Obermeyer, Z., Powers, B., Vogeli, C., Mullainathan, S.: Dissecting racial bias in an algorithm used to manage the health of populations. Science **366**(6464), 447–453 (2019)
38. Osgood, C.E.: Semantic differential technique in the comparative study of cultures. American anthropologist **66**(3), 171–200 (1964)
39. Ponzio, F., Deodato, G., Macii, E., Di Cataldo, S., Ficarra, E.: Exploiting "uncertain" deep networks for data cleaning in digital pathology. In: 2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI). pp. 1139–1143. IEEE (2020)
40. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021)
41. Rafailov, R., Sharma, A., Mitchell, E., Manning, C.D., Ermon, S., Finn, C.: Direct preference optimization: Your language model is secretly a reward model. Advances in Neural Information Processing Systems **36** (2024)
42. Rahimian, H., Mehrotra, S.: Distributionally robust optimization: A review. arXiv preprint arXiv:1908.05659 (2019)
43. Sagawa, S., Koh, P.W., Hashimoto, T., Liang, P.: Distributionally Robust Neural Networks. In: International Conference on Learning Representations (2020)
44. Sagawa, S., Koh, P.W., Lee, T., Gao, I., Xie, S.M., Shen, K., Kumar, A., Hu, W., Yasunaga, M., Marklund, H., et al.: Extending the wilds benchmark for unsupervised adaptation. arXiv preprint arXiv:2112.05090 (2021)
45. Sankaranarayanan, S., Hartvigsen, T., Oakden-Rayner, L., Ghassemi, M., Isola, P.: Real world relevance of generative counterfactual explanations. In: Workshop on Trustworthy and Socially Responsible Machine Learning, NeurIPS 2022
46. Schacter, D.L.: Implicit memory: History and current status. Journal of experimental psychology: learning, memory, and cognition **13**(3),  501 (1987)
47. Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C., Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis, C., Wortsman, M., et al.: Laion-5b: An open large-scale dataset for training next generation image-text models. Advances in Neural Information Processing Systems **35**, 25278–25294 (2022)
48. Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., Komatsuzaki, A.: Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. arXiv preprint arXiv:2111.02114 (2021)
49. Seo, S., Lee, J.Y., Han, B.: Unsupervised Learning of Debiased Representations with Pseudo-Attributes. In: IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 16742–16751 (2022)
50. Sohoni, N., Dunnmon, J., Angus, G., Gu, A., Ré, C.: No Subclass Left Behind: Fine-Grained Robustness in Coarse-Grained Classification Problems. Advances in Neural Information Processing Systems **33**, 19339–19352 (2020)
51. Stephan, W.G., Stephan, C.W.: Intergroup anxiety. Journal of social issues **41**(3), 157–175 (1985)
52. Thomee, B., Shamma, D.A., Friedland, G., Elizalde, B., Ni, K., Poland, D., Borth, D., Li, L.J.: Yfcc100m: The new data in multimedia research. Communications of the ACM **59**(2), 64–73 (2016)

53. Thompson, B.: Effect sizes, confidence intervals, and confidence intervals for effect sizes. Psychology in the Schools **44**(5), 423–432 (2007)
54. Vieriu, R.L., Tulyakov, S., Semeniuta, S., Sangineto, E., Sebe, N.: Facial Expression Recognition under a Wide Range of Head Poses. In: 2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG). vol. 1, pp. 1–7. IEEE (2015)
55. Xu, H., Xie, S., Tan, X.E., Huang, P.Y., Howes, R., Sharma, V., Li, S.W., Ghosh, G., Zettlemoyer, L., Feichtenhofer, C.: Demystifying clip data. arXiv preprint arXiv:2309.16671 (2023)
56. Yang, J.C., Korecki, M., Dailisan, D., Hausladen, C.I., Helbing, D.: Llm voting: Human choices and ai collective decision making. arXiv preprint arXiv:2402.01766 (2024)
57. Zhang, M., Sohoni, N.S., Zhang, H.R., Finn, C., Ré, C.: Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations. arXiv preprint arXiv:2203.01517 (2022)