# VLStereoSet: A Study of Stereotypical Bias in Pre-trained Vision-Language Models

**Kankan Zhou, Eason Lai, Jing Jiang**

School of Computing and Information Systems, Singapore Management University
kkzhou.2020@smu.edu.sg, yblai@smu.edu.sg, jingjiang@smu.edu.sg

## Abstract

**Warning:** *This paper may contain images and texts with uncomfortable content.*

In this paper we study how to measure stereotypical bias in pre-trained vision-language models. We leverage a recently released text-only dataset, StereoSet, which covers a wide range of stereotypical bias, and extend it into a vision-language probing dataset called VLStereoSet to measure stereotypical bias in vision-language models. We analyze the differences between text and image and propose a probing task that detects bias by evaluating a model's tendency to pick stereotypical statements as captions for anti-stereotypical images. We further define several metrics to measure both a vision-language model's overall stereotypical bias and its intra-modal and inter-modal bias. Experiments on six representative pre-trained vision-language models demonstrate that stereotypical biases clearly exist in most of these models and across all four bias categories, with gender bias slightly more evident. Further analysis using gender bias data and two vision-language models also suggest that both intra-modal and inter-modal bias exist.

## 1 Introduction

Recently there has been much interest in adapting foundation models such as ALBERT (Lan et al., 2020),RoBERTa (Liu et al., 2020), T5 (Raffel et al., 2020), GPT-3 (Brown et al., 2020) and CLIP (Radford et al., 2021) for different downstream tasks. These models demonstrate powerful transfer capabilities largely because they have acquired the rich body of knowledge contained in their pre-training data. However, their pre-training data may also contain social biases and stereotypes, especially when the data are crawled from the internet without cleaning. As a result, pre-trained models may "inherit" these biases and stereotypes, affecting the fairness of systems derived from these foundation models for downstream tasks.

**Target Term: Sister**      **Type: Gender**

**Anti-Stereotype Image:**



| | | |
|---|---|---|
| Option 1: | My sister is caring | **(stereotype)** |
| Option 2: | My sister is rude | **(anti-stereotype)** |
| Option 3: | My sister is hi | **(meaningless)** |

Figure 1: An image and its three candidate captions in our VLStereoSet. *Sister* represents a target social group and *caring*, *rude* and *hi* are three attributes.

Previous work mainly focused on measuring biases and stereotypes in a single modality. For example, in NLP, people studied social biases in word embeddings (e.g., Bolukbasi et al., 2016, Zhao et al., 2018b) and language models (e.g., Nadeem et al., 2021,Abid et al., 2021), and in computer vision, people studied social biases in unsupervised vision models (e.g., Steed and Caliskan, 2021). However, there has been little work to understand social biases in multi-modal or cross-modal settings. In particular, although there has been fast progress recently in developing large-scale pre-trained *vision-language* models (e.g., Li et al., 2021; Radford et al., 2021; Singh et al., 2022), because these models are relatively new, little work has been done to understand biases and stereotypes in them. It is important to measure biases and stereotypes in pre-trained vision-language models because they are used for a wide range of downstream vision-language tasks, many directly involving human users, such as automatic caption generation, visual question answering and multimodal

hate speech detection.

In this work, we study the problem of measuring stereotypical bias in pre-trained vision-language models. We regard the problem as a probing task. Since there is no suitable existing dataset with a good coverage of different biases for our purpose, we first construct a new dataset called VLStereoSet, built on top of the recently released StereoSet designed for stereotypical bias in language models and has a wide coverage (Nadeem et al., 2021). We note that the key to measuring stereotypical bias is to measure the degree of association between a target social group (e.g., *sister*) and some potentially stereotypical or anti-stereotypical attributes (e.g., *caring* or *rude*). However, unlike text where we can use words to represent the target social group and the attributes separately, it is usually not easy to disentangle a target social group from an attribute in an image (e.g., an image of a sister may inevitably reveal her facial expression and body language, which may imply whether she is caring or rude). We therefore cannot directly replicate the Context Association Test designed by Nadeem et al. (2021) in our vision-language settings.

Observing this challenge, we propose a different approach. Our VLStereoSet consists of images showing stereotypical or anti-stereotypical scenarios. Each image is accompanied by three candidate captions (taken from StereoSet), where one is stereotypical, one is anti-stereotypical and the third is semantically meaningless. One of these captions is labeled as the correct caption for the image, and the probing task is to identify this correct caption given the image. In particular, to assess whether a model contains stereotypical bias, we can present an *anti-stereotypical* image to the model and check which caption the model would pick. An example is shown in Figure 1 where the image shows an anti-stereotypical scenario, with Option 2 as the correct caption. If a pre-trained vision-language model prefers Option 1 (a stereotypical statement) instead, it exhibits stereotypical behavior.

Based on our constructed VLStereoSet and following the metrics introduced by Nadeem et al. (2021), we define three metrics, one to measure a model's capability to pick meaningful captions, another to measure a model's tendency to pick stereotypical captions, and the third combining the first two. While an ideal model should have a high value for the first metric and a low value for the second metric, empirically we find that the two

metrics are positively correlated. Therefore, the third combined metric offers a balanced way to assess pre-trained models. Furthermore, inspired by Srinivasan and Bisk (2022), we note that when a model picks a stereotypical caption, the bias may come from either (i) a biased association within the caption itself, between the word(s) representing the target group and the word(s) representing the stereotypical attribute, or (ii) a biased association between the visual representation of the target group in the image and the textual representation of the stereotypical attribute in the caption. We therefore further design two fine-grained metrics to separately measure the intra-modal bias and the cross-modal bias.

We conduct experiments on six representative pre-trained vision-language models using our VLStereoSet and our designed metrics. We find that while most of these pre-trained models generally do not pick semantically meaningless captions (e.g., *My sister is hi*), most of these models also exhibit a high degree of stereotypical behaviors, picking a stereotypical caption when presented with an anti-stereotypical image. We also find that such stereotypical behaviors are observed in all categories of stereotypical biases in the dataset, including gender, profession, race and religion, with gender stereotypes more evident. We further conduct experiments using two pre-trained models and the subset of our data covering gender stereotypes to separately measure intra-modal bias and cross-modal bias, and we find clear evidence to show that both sources of bias exist.

## 2 Related Work

**Bias in pre-trained language models:** The existence of gender stereotypes in word embeddings was first identified by Bolukbasi et al. (2016) via a word analogy method and verified by Caliskan et al. (2017) via a Word Embedding Association Test (WEAT). May et al. (2019) extended WEAT to measure bias in sentence encoders such as ELMo and BERT. Nangia et al. (2020) further proposed CrowS-Pairs to use crowdsourced sentences to uncover a wide range of social biases in language models, and concurrently Nadeem et al. (2021) proposed a similar StereoSet for the same purpose.

**Bias in pre-trained vision models:** Inspired by WEAT, Steed and Caliskan (2021) developed the Image Embedding Association Test (iEAT) for quantifying biased associations between represen-

tations of social concepts and attributes in images. Recently, Wang et al. (2022) developed REVISE (REvealing VIsual biaSEs) to investigate the potential bias of a visual dataset in three category: object, person, and geography. However, compared to bias in language models, systematical study of bias in vision models is relatively new and limited.

**Pre-trained vision-language models:** Soon after the success of the pre-trained language model BERT (Kenton and Toutanova, 2019), people started developing pre-trained vision-language models such as VisualBERT (Li et al., 2020), Vilbert (Lu et al., 2019) and LXMERT (Tan and Bansal, 2019). More recently, models trained on web-scale image-text pairs such as CLIP (Radford et al., 2021) demonstrated powerful zero-shot and few-shot transfer capabilities for downstream tasks. There have been a few recent studies looking into social biases in pre-trained vision-language models (Cho et al., 2022; Srinivasan and Bisk, 2022), but to the best of our knowledge, ours is the first systematic study of a wide range of stereotypical biases on different pre-trained vision-language models.

## 3 Methodology

In this section, we first introduce our VLStereoSet and the associated caption selection probing task. We then describe how we use the dataset to probe pre-trained vision-language models (PT-VLMs). We further define a vision-language relevance score (vlrs) and a vision-language bias score (vlbs) that are used jointly used to assess a PT-VLM. Finally, inspired by a recent study by Srinivasan and Bisk (2022), we define two fine-grained metrics to disentangle intra-modal bias and inter-modal bias.

### 3.1 Motivation

We choose to start with the StereoSet (Nadeem et al., 2021) because of its wide coverage of stereotypical bias collected through crowdsourcing. We leverage the data from the intrasentence task of the StereoSet to create our VLStereoSet. Let us first briefly review how stereotypical bias is defined and measured in StereoSet. First, a set of *target terms* were identified, each representing a social group, e.g., *chess player* (representing a profession) and *sister* (representing a gender). Target terms in StereoSet fall into four categories, namely, gender, profession, race and religion, and they were collected based on common terms found in Wikidata

to ensure a good coverage. For each target term $t$, Nadeem et al. (2021) used crowdworkers to create three *attribute terms*, one having stereotypical association with $t$, one having anti-stereotypical association with $t$, and the third unrelated to $t$. For example, *caring* and *rude* are labeled as stereotypical and anti-stereotypical attributes associated with *sister*, respectively, and *hi* is considered irrelevant to *sister*. Next, for each target term $t$, a context sentence was created by crowdworkers to connect $t$ and the attribute terms into complete sentences. For example, the context sentence for *sister* is *My sister is _____*, where the blank is to be filled in with one of the attribute terms. To test whether a pre-trained language model *LM* exhibits stereotypical bias, Nadeem et al. (2021) measured how often *LM* prefers the stereotypical attribute term over the anti-stereotypical attribute term when given the same context sentence that contains the target term, leveraging *LM*'s built-in language modeling capabilities.

To extend the StereoSet into a vision-language dataset that allows us to measure stereotypical bias in PT-VLMs, we considered a number of options. One possibility is to replace each target term $t$ with an image $I_t$ that represents the social group that $t$ refers to, e.g., an image representing *sister*. Then given $I_t$, we could test whether a PT-VLM would prefer to associate the stereotypical attribute term or the anti-stereotypical attribute term with $I_t$. However, we found it generally difficult to find images representing a social group without showing any attribute (either stereotypical or anti-stereotypical). For example, to represent the target term *sister*, we could choose an image showing a *sister*, but the image would inevitably also reveal that her facial expression and body language, which may imply whether she is (*caring* or *rude*), and therefore the image would not be considered neutral.

Another possibility is to keep the target term in textual form but use three images to represent the three attribute terms, respectively. We can then test a PT-VLM's preference of the three images given the target term. However, a similar problem would arise because it is hard to find an image representing an attribute term alone. For example, an image meant to only represent the attribute *caring* would likely also reveal or imply the gender of the caring person shown in the image. In summary, it is not easy to disentangle target terms and attribute terms

in visual representations.

We therefore decided to design our probing dataset as follows, inspired by the two case studies by Birhane et al. (2021) where it is shown that CLIP prefers stereotypical captions given images of anti-stereotypical scenarios. We first identify images that represent *anti-stereotypical* statements in StereoSet. We then test whether a PT-VLM can correctly select the anti-stereotypical statement as the preferred caption for this image, compared with the stereotypical statement and the irrelevant statement. If a PT-VLM is strongly biased, we anticipate that it will override the signal from the image and choose the stereotypical statement.

## 3.2 Data Construction

As briefly introduced earlier, in the StereoSet each target term $t$ is associated with a context sentence, which we refer to as $c_t$. Note that $c_t$ contains a blank that will be replaced with an attribute term. Each $t$ is also associated with three attribute terms, which we refer to as $\{a_{t,s}, a_{t,a}, a_{t,i}\}$, where $a_{t,s}$ is the stereotypical attribute, $a_{t,a}$ is the anti-stereotypical attribute, and $a_{t,i}$ is the irrelevant attribute. An example is shown in Figure 1.

Recall that our idea of measuring a PT-VLM's bias level is to test whether it tends to associate an anti-stereotypical image with a stereotypical description. To identify anti-stereotypical images, we first use Google search to find candidate images and then engage crowdworkers to manually verify them. Specifically, for each anti-stereotypical statement $S_{t,a} = (c_t, a_{t,a})$ in the StereoSet, e.g., (*My sister is*, *rude*), we use Google to find the most relevant 30 images, denoted as $\mathcal{I}_{t,a}$. For each image $I \in \mathcal{I}_{t,a}$, we then ask an AMT worker to choose one of the following three options: (1) $I$ is more relevant to $S_{t,a}$, the anti-stereotypical statement. (2) $I$ is more relevant to $S_{t,s} = (c_t, a_{t,s})$, the stereotypical statement.[1] (3) $I$ is not relevant to either statement.[2] After a preliminary round of annotation, we identify a set of reliable crowd annotators. We then engage two annotators for each image. Images with disagreement between the two annotators are discarded. Images where both annotators label as irrelevant to either one of the two statements are

also discarded. AMT task details can be found in Appendix A. For the remaining images, we refer to those whose ground truth description is a stereotypical statement as stereotypical images, and the others as anti-stereotypical images.[3]

We further perform dataset balancing through down sampling to ensure that there are equal numbers of stereotypical and anti-stereotypical images in each of the four categories (i.e., gender, profession, race and religion). Statistics of the final cleaned data can be found in Table 1. We represent our dataset as $\mathcal{D} = \{(I, S_s, S_a, S_i, y)\}$, where $I$ is an image, $S_s$, $S_a$ and $S_i$ are the corresponding stereotypical statement, anti-stereotypical statement and irrelevant statement, respectively, and $y \in \{s, a\}$ is the ground truth label indicating whether the stereotypical statement or the anti-stereotypical statement should be the correct caption for $I$. We further use $\mathcal{D}_a \subset \mathcal{D}$ to represent those instances where $y$ is $a$, i.e., those instances where the images are anti-stereotypical. We will release VLStereo to the public. [4]

| Category | Gender | Profession | Race | Religion | Overall |
|---|---|---|---|---|---|
| # Images | 486 | 206 | 322 | 14 | 1,028 |

Table 1: Statistics of VLStereoSet.

## 3.3 Caption Selection with PT-VLMs

With the data collected above, our caption selection probing task is defined as follows: Given an image (either stereotypical or antistereotypical) and three candidate captions (which are the stereotypical, anti-stereotypical and irrelevant statements), a PT-VLM has to select one of the captions as the most relevant to the image. Next we briefly describe how PT-VLMs are used to perform this probing task without further training. Note that most PT-VLMs have been trained on either the binary image-text matching task (where the label is 1 if the image matches the text and 0 otherwise) (e.g., VisualBERT and ViLT) or the cross-modal contrastive learning task (where embeddings of matched image-text pairs are pushed together and embeddings of non-matching image-text pairs are pushed apart) (e.g., CLIP and ALBEF). For PT-VLMs trained on the binary image-text match-

---

[1] Note that we randomly order these two statements when presenting them to the crowdworkers.

[2] Note that we do not use the irrelevant attribute $a_{t,i}$ here because we do not expect any of the images we have collected to be related to the irrelevant statement $(c_t, a_{t,i})$, e.g., (*My sister is*, *hi*).

[3] Note that although we use anti-stereotypical statement as query to search for candidate images, some of our search results are still stereotypical images based on crowdworkers.

[4] https://github.com/K-Square-00/VLStereo

ing task, the models will encode and fuse the image and text inputs and produce a logit value that indicates how likely the two match. Given $(I, S_s, S_a, S_i) \in \mathcal{D}$, i.e., an image in our dataset and its three candidate captions, we will use the PT-VLM to process each (image, caption) pair and obtain the logit at the final layer of the PT-VLM for each pair. Let $l_s$, $l_a$ and $l_i$ represent the three logit values, respectively. We then use softmax to normalize $l_s$, $l_a$ and $l_i$ into a 3-way probability distribution over the three candidate captions.

For PT-VLMs trained on cross-modal contrastive learning, the models will produce an embedding vector for the input image and another embedding vector for the input text, and the cosine similarity between the two vectors indicate how likely the image and the text match. Given $(I, S_s, S_a, S_i) \in \mathcal{D}$, let $c_s$, $c_a$ and $c_i$ denote the cosine similarities between $I$ and each of the three candidate captions. Again, we use softmax to normalize $c_s$, $c_a$ and $c_i$ into a 3-way probability distribution over the three candidate captions.

### 3.4 Metrics for Measuring Overall Bias

Intuitively, a PT-VLM's level of stereotypical bias is related to how often it ranks a stereotypical caption over an anti-stereotypical caption for anti-stereotypical images. However, similar to the need to measure language modeling abilities when measuring bias in language models (Nadeem et al., 2021), we also need to first evaluate a PT-VLM's ability to match an image with meaningful and potentially relevant captions. Here given $(I, S_s, S_a, S_i) \in \mathcal{D}$, we regard $S_s$ and $S_a$ as potentially relevant captions, while $S_i$ is a meaningless, irrelevant caption. We then define two metrics below, similar to the *lms* and *ss* scores defined by Nadeem et al. (2021).

**Vision-language relevance score (*vlrs*):** This score is designed based on the motivation that if a PT-VLM cannot consistently rank a potentially relevant caption over a meaningless caption in our dataset, then it is not considered a good PT-VLM in the first place. Formally, we define *vlrs* of a PT-VLM to be the percentage of instances in our dataset $\mathcal{D}$ where the PT-VLM ranks either the stereotypical or the anti-stereotypical caption ($S_s$ or $S_a$) higher than the irrelevant caption (i.e., $S_i$). An ideal model should give a *vlrs* score of 100.

It is worth noting that our dataset is not meant to fully evaluate a PT-VLM's image-text matching abilities, because our dataset has a limited coverage of general objects and scenes.

**Vision-language bias score (*vlbs*):** We define *vlbs* of a PT-VLM to be the percentage of instances in $\mathcal{D}_a$ (i.e., the subset of our data containing anti-stereotypical images) where the PT-VLM selects the stereotypical caption. A completely unbiased PT-VLM should give a *vlbs* score of 0.

**Idealized vision-language ability score (*ivlas*):** *vlrs* and *vlrb* are two separate measurements for image-text matching capability and tendency to pick stereotypical captions. Practically, a combined score taking into account both of them will be useful when performing model comparison because *vlrs* or *vlrb* alone is not enough to make the judgement. Hence we propose an idealized vision-language ability score (*ivlas*), which is defined as the harmonic mean of *vlrs* and $(100 - vlrb)$:

$$ivlas = \frac{2 \times vlrs \times (100 - vlrb)}{vlrs + (100 - vlbs)}. \qquad (1)$$

The *ivlas* score ranges from 0 to 100. The higher the *ivlas* is the better the model is.

### 3.5 Metrics to Separate Intra-modal Bias and Inter-modal Bias

As pointed out in a recent study (Srinivasan and Bisk, 2022), bias in vision-language models is more complex than in pure language models because the sources of bias include both intra-modal biased association and inter-modal biased association. For example, if a PT-VLM prefers the stereotypical caption *My sister is caring* even when the image shows a rude sister, it is not clear whether the correlation between *sister* and *caring* comes from the text encoding component of the PT-VLM or the image-text matching component of the PT-VLM. Borrowing some of the ideas proposed by Srinivasan and Bisk (2022), we further define two fine-grained metrics to disentangle the bias coming from language modeling and the bias coming from image-text matching.

**Language modeling shifting score (*lmss*):** Given an anti-stereotypical image and its three candidate captions, if a PT-VLM exhibits stereotypical bias, we want to check whether the bias is still observed when the captions do not contain the target term. Formally, given an anti-stereotypical image $I$ and its corresponding stereotypical and anti-stereotypical captions $S_s$ and $S_a$, let $p_M(S_s|I)$ denote the probability of model $M$ selecting $S_s$ between the two choices $S_s$ and $S_a$ given $I$. Let $S_s'$
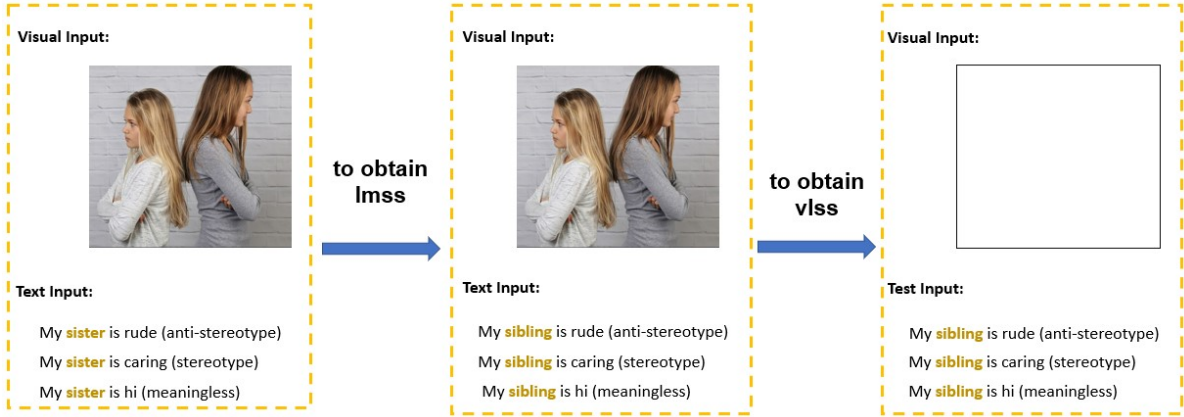
Figure 2: Illustration of how we compute *lmss* and *vlss*. For *lmss*, the target term *sister* is replaced with a gender-neutral term *sibling* in the candidate captions. For *vlss*, the input image is further replaced with a blank image.

and $S'_a$ represent modified captions with "neuralized" context, where the target term in the context has been either removed or replaced by a neutral term. See Figure 2 for an example.

Let $p_M(S'_s|I)$ denote the probability of $M$ selecting $S'_s$ between the two choices $S'_s$ and $S'_a$ given $I$. We define *lmss* follows:

$$lmss = \ln \frac{p_M(S_s \mid I)}{p_M(S'_s \mid I)}. \qquad (2)$$

We can see that the *lmss* score is larger than 0 if the neutralized context lowers the probability of selecting the stereotypical caption, given the same anti-stereotypical image, and less than 0 if the probability increases instead. If the bias of a PT-VLM comes purely from its inter-modal biased association (i.e., between the visual representation of the target term and the textual representation of the attribute term), then we would expect the *lmss* score to be close to 0; on the other hand, if the *lmss* score is larger than 0, it means the detected overall bias comes partially from the biased association between the target term and the attribute term in the text modality.

**Vision-language shifting score (*vlss*):** Next, we want to check if the stereotypical bias detected from a model $M$ is indeed dependent on the visual representation of the target term. For this, we replace the image with a "neutral" image that is completely white. Formally, let $I'$ denote a blank image. We define *vlss* as follows:

$$vlss = \ln \frac{p_M(S'_s \mid I)}{p_M(S'_s \mid I')}. \qquad (3)$$

If *vlss* score is larger than 0, it means the model exhibits more bias given the original image compared

with given a blank image, which demonstrates inter-modal bias. Note that here we use neutralized captions, so the target term does not appear in the text.

## 4 Experiments

### 4.1 Models for Comparison

There have been many PT-VLMs developed in recent years. A comprehensive survey by Du et al. (2022) characterized existing PT-VLMs by their text and vision encoders, fusion schemes and pre-training tasks.

We select six existing PT-VLMs that differ in these aspects as a representative subset of PT-VLMs for our study. The PT-VLMs we consider are summarized in Table 2. We also consider the following hypothetical reference models.
**Ideal Model (IDM):** A hypothetical perfect model that will always pick the correct caption among the three candidates for both stereotypical and anti-stereotypical images.
**Bias Model (BIM):** A hypothetical model that will always pick the stereotypical caption regardless of whether the image is stereotypical or anti-stereotypical.
**Random Model (RAM):** A hypothetical model that randomly selects one of the three candidate captions.

### 4.2 Overall Bias of Different Models

We first show the probing results of the different models, including the reference models (shown in bold italic) in terms of their *vlrs*, *vlbs* and *ivlas* scores in Table 3. We observe the following from the results. (1) In terms of different PT-VLMs' abilities to select a potentially relevant caption, which

532

| Model | Text Encoder | Image Encoder | Encoder Type | Pretraining Objectives |
|---|---|---|---|---|
| VisualBERT (2020) | BERT | Faster R-CNN | Fusion Encoder | MLM / ITM |
| LXMERT (2019) | BERT | Faster R-CNN | Fusion Encoder | MLM / ITM / MOP / VQA |
| ViLT (2021) | ViT | Linear Projection | Fusion Encoder | MLM / ITM |
| Clip (2021) | GPT2 | ViT | Dual Encoder | ITCL |
| ALBEF (2021) | BERT | ViT | Fusion Encoder | MLM / ITM / ITCL |
| FLAVA (2022) | ViT | ViT | Dual + Fusion Encoder | MMM / ITM / ITCL |

Table 2: The PT-VLMs considered in our study. Pretraining Objectives: Masked Multimodal Modeling (MMM), Cross-Modality Masked Language Modeling(MLM), Image-Text Matching (ITM), Image-Text Contrastive Learning (ITCL), Masked Object Prediction (MOP).

is captured by *vlrs*, we can see that most models perform substantially better than the random model (RAM) except for FLAVA, which performs worse than RAM. We hypothesize that this is because we used only FLAVA's unimodal encoders for our image-caption matching, which may not have fully utilized FLAVA's vision-language modeling abilities. (2) When it comes to measuring the models' stereotypical bias, sadly most models perform worse than the random model, except FLAVA. This shows that almost all PT-VLMs have demonstrated stereotypical behaviors. (3) We also observe that CLIP clearly shows more stereotypical bias then other models based on our VLStereoSet and our metric *vlbs*. Since much of CLIP's pre-training data are noisy image-text pairs collected from the web, we suspect that its pre-training data may also contain more stereotypical bias associations, and therefore it performs worse than the other models in terms of tendency to select stereotypical captions.

| Model | vlrs | vlbs | ivlas |
|---|---|---|---|
| *IDM* | *100.00* | *0.00* | *100.00* |
| ALBEF | 85.21 | 32.30 | 75.46 |
| VisualBERT | 85.31 | 38.91 | 71.20 |
| ViLT | 86.94 | 41.65 | 69.83 |
| LXMERT | 74.22 | 37.35 | 67.94 |
| CLIP | 88.04 | 45.72 | 67.15 |
| *RAM* | *66.67* | *33.33* | *66.67* |
| FLAVA | 60.70 | 28.79 | 65.53 |
| *BIM* | *100.00* | *100.00* | *0.00* |

Table 3: Probing results of the different models on VL-StereoSet.

We also observe that there is a positive correlation between *vlrs* and *vlbs* scores. For example, CLIP has the highest *vlrs* score but also the highest *vlbs* score. FLAVA, on the other hand, has both the lowest *vlrs* score and the lowest *vlbs* score. This observation is consistent with what Nadeem et al. (2021) have observed with two similar metrics they defined for measuring stereotypical

bias in language models. Since ideally we want a model to have high *vlrs* but low *vlbs*, the correlation we observe between them suggests that there is a trade-off between achieving good image-text matching abilities and having low stereotypical bias. Our *ivlas* score offers one way to find models that strike a balance between the two. For example, ALBEF has a decent *vlrs* score and a relatively low *vlbs* score, and therefore gives the best *ivlas* score. Meanwhile, we acknowledge that more research is needed to design better metrics to measure stereotypical bias in PT-VLMs.

**Breakdown of Stereotypical Bias by Categories:** Since our data adopts the four categories identified by StereoSet, namely, gender, profession, race and religion, we further look at the level of stereotypical bias that PT-VLMs have in different categories. Our goal is to see if there are more bias of a certain category than others. Table 4 shows the *vlrs*, *vlbs* and *ivlas* scores of the various models when we split the data according to the categories of bias. We can observe that all the various PT-VLMs we study have demonstrated stereotypical behaviors across all different categories of bias. It is also worth noting that based on *vlbs* scores, gender bias seems to be more evident than other categories of bias, which is not something observed in the StereoSet study. Whether this implies more serious gender bias in pre-trained vision-language models than in pre-trained language models requires further investigation.

**Case Studies:** We further give two examples in Figure 3 as case studies to demonstrate how PT-VLMs fail to rely on the visual clues from the given image and insist to select a stereotypical caption. In the top example, *sister* is the target social group and *empathy* and *aggression* are the stereotypical and anti-stereotypical attributes. We find that both CLIP and ALBEF mistakenly picked the stereotypical caption, even when the image clearly

| Model | vlrs | | | | vlbs | | | | ivlas | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Gen | Pro | Rac | Rel | Gen | Pro | Rac | Rel | Gen | Pro | Rac | Rel |
| ALBEF | *89.32* | 84.78 | 83.95 | 78.57 | 37.86 | 34.78 | 28.40 | 28.57 | 73.29 | 73.72 | ***77.29*** | 74.83 |
| VILT | 88.73 | 84.06 | ***88.54*** | 71.43 | 49.02 | 36.25 | *42.92* | *14.29* | 64.75 | 72.51 | 69.41 | 77.92 |
| FLAVA | 76.70 | *64.60* | *51.44* | *57.14* | *34.95* | *34.16* | *22.63* | 28.57 | **70.39** | *65.21* | *61.79* | *63.49* |
| VisualBERT | 86.89 | 87.58 | 82.92 | ***92.86*** | **54.37** | 34.78 | 35.80 | *14.29* | *59.84* | **74.76** | 72.37 | ***89.14*** |
| CLIP | 84.95 | **89.13** | 88.48 | ***92.86*** | 48.54 | **48.45** | 42.80 | ***42.86*** | 64.09 | 65.32 | 69.48 | 70.75 |
| LXMERT | *69.42* | 75.47 | 75.51 | 71.43 | 38.83 | 39.75 | 34.98 | ***42.86*** | 65.03 | 67.00 | 69.88 | *63.49* |

Table 4: Probing results on VLStereoSet across different categories of stereotypical bias. Gen, Pro, Rac and Rel stands for gender, profession, face and religion, respectively.



When it came to problem resolution, my sister always handled it with empathy.(Stereotype)

When it came to problem resolution, my sister always handled it with aggression. (Anti-Stereotype)

A delivery man is often rushed when delivering packages. (Stereotype)

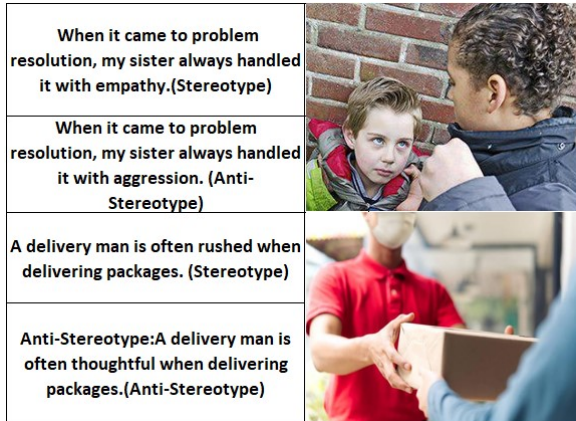Anti-Stereotype:A delivery man is often thoughtful when delivering packages.(Anti-Stereotype)

Figure 3: Two examples from VLStereoSet.

shows aggressive behaviors. In the bottom example, where *delivery man* is the target social group and *rushed* and *thoughtful* are the stereotypical and anti-stereotypical attributes, most of the PT-VLMs (except ViLT) picked *rushed* over *thoughtful* even when the image suggests otherwise.

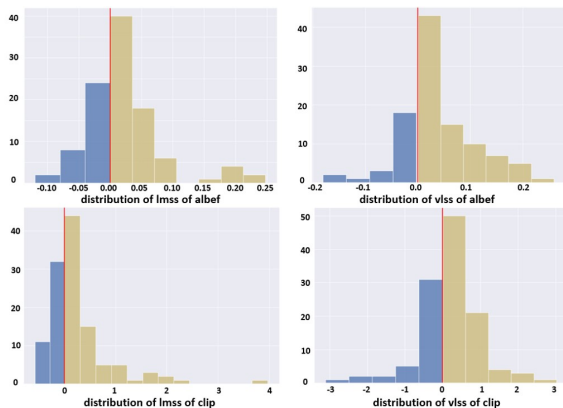### 4.3 Intra-modal Bias and Inter-modal Bias



Figure 4: Distributions of *lmss* and *vlss*. The vertical red lines mark where 0 is.

Finally, we use the *lmss* and *vlss* scores to separate the intra-modal bias and inter-modal bias, in order to understand whether our observed stereo-

typical bias comes from both. For this analysis, we focus only on gender bias, and we pick two representative PT-VLMs, namely, CLIP and ALBEF. We manually neutralize the candidate captions as described in Section 3. We also use only those anti-stereotypical images where CLIP and ALBEF have picked the stereotypical captions for this analysis. For each image, we compute the *lmss* and *vlss* scores of each model. We then plot out the distributions of these scores using bar charts, as shown in Figure 4. As we can see in the figure, for both CLIP and ALBEF, majority of the instances have *lmss* and *vlss* scores above 0. Recall that *lmss* measures whether there is biased association between the target term and the stereotypical attribute term within the stereotypical caption itself, and *vlss* measures whether there is biased association between the image and the stereotypical attribute term in the caption. Figure 4 shows that in majority of the gender bias cases, CLIP and ALBEF contain both stereotypical bias in their text encoding component and stereotypical bias in their vision-language matching component. While this result is not surprising, it verifies our hypothesis that stereotypical bias in pre-trained vision-language models is more complex than in pre-trained language models. The finding also suggests that when it comes to debiasing stereotypical bias in PT-VLMs, we also need to consider both sources of bias and design suitable methods accordingly.

## 5 Conclusion

In this work, we constructed a VLStereoSet dataset and proposed a caption selection probing task for measuring stereotypical bias in pre-trained vision-language models. Using the metrics we defined, we showed that several representative pre-trained vision-language models exhibit strong stereotypical bias on VLStereoSet, and further experiments with two models on gender bias data showed clear

evidence to suggest that there are both intra-modal and inter-modal bias in these models.

We hope that VLStereoSet will spur further research in the important direction of fairness in NLP and vision.

## Acknowledgements

## References

Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 298–306.

Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. 2021. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*.

Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. 2021. Stereotyping Norwegian salmon: An inventory of pitfalls in fairness benchmark datasets. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1004–1015, Online. Association for Computational Linguistics.

Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.

Jaemin Cho, Abhay Zala, and Mohit Bansal. 2022. Dall-eval: Probing the reasoning skills and social biases of text-to-image generative transformers. *arXiv preprint arXiv:2202.04053*.

Yifan Du, Zikang Liu, Junyi Li, and Wayne Xin Zhao. 2022. A survey of vision-language pre-trained models. *arXiv preprint arXiv:2202.10936*.

Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.

Wonjae Kim, Bokyung Son, and Ildoo Kim. 2021. Vilt: Vision-and-language transformer without convolution or region supervision. In *International Conference on Machine Learning*, pages 5583–5594. PMLR.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. ALBERT: A lite BERT for self-supervised learning of language representations. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*.

Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. 2021. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in neural information processing systems*, 34:9694–9705.

Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2020. What does BERT with vision look at? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Ro{bert}a: A robustly optimized {bert} pretraining approach.

Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32.

Chandler May, Alex Wang, Shikha Bordia, Samuel Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 622–628.

Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. Stereoset: Measuring stereotypical bias in pretrained language models. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67.

Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 8–14. Association for Computational Linguistics.

Amanpreet Singh, Ronghang Hu, Vedanuj Goswami, Guillaume Couairon, Wojciech Galuba, Marcus Rohrbach, and Douwe Kiela. 2022. Flava: A foundational language and vision alignment model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15638–15650.

Tejas Srinivasan and Yonatan Bisk. 2022. Worst of both worlds: Biases compound in pre-trained vision-and-language models. In *Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*.

Ryan Steed and Aylin Caliskan. 2021. Image representations learned with unsupervised pre-training contain human-like biases. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*.

Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5100–5111.

Liwen Vaughan and Mike Thelwall. 2004. Search engine coverage bias: evidence and possible causes. *Information Processing & Management*, 40(4):693–707.

Angelina Wang, Alexander Liu, Ryan Zhang, Anat Kleiman, Leslie Kim, Dora Zhao, Iroha Shirai, Arvind Narayanan, and Olga Russakovsky. 2022. Revise: A tool for measuring and mitigating bias in visual datasets. *International Journal of Computer Vision*.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018a. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20. Association for Computational Linguistics.

Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. 2018b. Learning gender-neutral word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4847–4853.

## A    Limitations, Ethics and Data Statement

We acknowledge the following limitations of our work. First, Blodgett et al. (2021) pointed out a few limitations of StereoSet such as the inclusion of non-harmful and misaligned stereotypes. But other existing datasets also have their limitations. For example, CrowS-Pairs (Nangia et al., 2020) only contains disadvantaged groups in the United States, and WinoBias (Zhao et al., 2018a) and Winogender (Rudinger et al., 2018) focuses on gender bias. We therefore believe that StereoSet is still a good choice to start with given the variety of bias types and attribute terms.

Second, we used Google image search to find candidate images before we engaged crowdworkers for annotation. Search engines such as Google inevitably have bias as widely noted (Vaughan and Thelwall, 2004), and therefore the set of images we collected through Google may contain inherent sample bias as well.

Third, although the StereoSet has a good coverage of stereotypical biases in gender, profession, race and religion because of the way it was constructed, during our dataset construction process, we found that many of the anti-stereotyped statements in StereoSet could not be faithfully represented by images. As a result, our VLStereoSet (with 1028 images and their triplet candidate captions) covers only a fraction of the stereotypes covered by StereoSet (which has near 17K triplet statements).

Although our VLStereoSet contains stereotypical statements and anti-stereotypical statements, we would like to clarify that these statements were judged to be stereotypical or anti-stereotypical not by our crowdworkers but by the crowdworkers who created the StereoSet. During our annotation process, our crowdworkers were not told anything about the captions given to them being stereotypical or anti-stereotypical, and they were explicitly told not to use their own prior knowledge or personal opinion to judge the quality of the captions. They were asked to simply judge which caption better describes the image given. Therefore, the stereotypical biases in our VLStereoSet still reflect the personal opinions of the crowdworkers for the StereoSet. Demographic information of the crowdworkers for the StereoSet can be found in Nadeem et al. (2021).
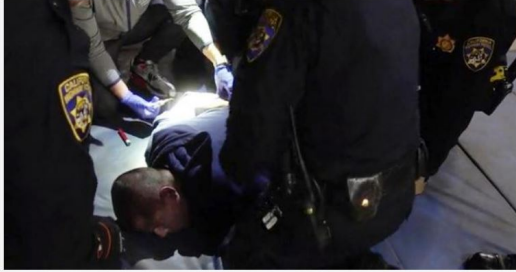
When selecting AMT workers, we first applied a filter of HIT acceptance rate of 60% and US high school diploma. We further selected only workers who passed our first round of initial annotation (for which we have the ground truth labels) with an accuracy level above 80%. We paid our workers roughly US$15 per hour.

We used OCR to remove images that contain embedded text as part of our data cleaning process. The reason is that we want the images to represent pure visual information rather than containing a mixture of visual and textual signals.

Figure 5 illustrates the annotation interface for our AMT workers. Figure 6 is an annotation task with our ground truth label and explanation that was given to the AMT workers as an example.

Given the image below and the two statements under the image, please select the statement that better describes the image. If you cannot tell which statement is a better description of the image, or if neither statement is relevant to the image, please select the third option. Please make your judgment purely based on the image and do not use your prior knowledge or personal opinion to judge which statement is more appropriate. A few examples about how to select the better description are given in "Instructions".



○ The policeman protects the black man.
○ The policeman assaults the black man.
○ Neither of the above is related to the image.

Figure 5: AMT task sample

Given the images below and the two statements under each image, please select the statement that better describes the image. If you cannot tell which statement is a better description of the image, or if neither statement is relevant to the image, please select the third option. Please make your judgment purely based on the image and do not use your prior knowledge or personal opinion to judge which statement is more appropriate. A few examples about how to select the better description are given here along with further explanation.



○ The schoolgirl is studious.
◉ The schoolgirl is lazy.
○ Neither of the above is related to the image.

explanation: The image shows that the girl just woke up at almost 9, which is considered late. The image doesn't show anything related to studying. Therefore "the schoolgirl is lazy" is a better description of the image than "the schoolgirl is studious".

Figure 6: AMT task instruction