# SocialCounterfactuals: Probing and Mitigating Intersectional Social Biases in Vision-Language Models with Counterfactual Examples

## Supplementary Material

## 7. Additional Analysis of SocialCounterfactuals

### 7.1. Examples of Counterfactual Sets

Figure 14, 15 and 16 provides additional examples of counterfactual sets generated by our approach.

### 7.2. Error Analysis

To investigate the different failure modes for the generated images in this study, we conducted a human evaluation of counterfactual sets for gender-race bias in occupations. We sampled 1200 images (100 for each gender-race combination) prior to our third stage of filtering (CLIP Attribute Detectability Filtering) and then annotated them into 5 categories. Results are shown in Table 6. We found that 90.8% of the images were correctly generated in terms of occupation, gender and race. In 5.2% of the samples, the gender could not be identified. This was typically due to subjects looking backwards or facing the wrong direction.

*Failure to generate female subjects* was the second most frequent error with 2.2%, followed by *subject completely out of frame/focus*. The least common error was *Failure to generate male subject* with a frequency of only 0.8%. No failures related to race were observed. Sampled images illustrating each of the different modes of failures are displayed in Table 7.

| Error Category | % present in sample |
|---|---|
| Good | 90.8% |
| Cannot discern gender | 5.2% |
| Failure to generate female subject | 2.2% |
| Subject completely out of frame/focus | 1.0% |
| Failure to generate male subject | 0.8% |

Table 6. Error analysis for 1200 random samples focused on gender and race.

An extended breakdown by race is shown in Table 8. We observe that across White, Black, Indian and Asian races, failure to generate either female or male subjects is approximately the same. For Latino images, there is a higher proportion of failures to generate female subjects; for Middle Eastern (M.E.) subjects, the proportion is even more pronounced, with a 12% failure to generate female subjects. In contrast, no generation failures were observed for Middle Eastern male subjects.

To further study the impact of our filtering method on image quality, we measure the same generation failure percentages after applying CLIP Attribute Detectability Filtering.

This final stage of filtering increases the percentage of *Good* images to 99% and 96% for male and female (respectively). Since the proportion of male and female are the same, the overall accuracy increases from 90.8% to 97.5% when CLIP Attribute Detectability Filtering is used.

### 7.3. Analysis of Occupation Detectability

We manually evaluated images depicting all of the occupations listed in Table 18 to identify those which are hard to recognize even for humans. In randomly selected subsets, we observed that occupations such as accountant, salesperson, economist, broker, web developer, attorney, banker, mathematician, and engineer produce very similar images depicting subjects wearing formal attire. Although the images accurately depict the social attributes (e.g., race, gender), the subject can be difficult to distinguish from other similar occupations. Because we evaluate skewness separately for each occupation, this similarity across occupations should not have a significant impact on the validity of bias estimation.

We also observed that it can be difficult to distinguish subjects across the various the medical profession, (e.g., pediatrician, audiologist) since they are depicted by Stable Diffusion as generic doctors wearing coats or scrubs. Similar depictions were also seen for education-related professions (e.g., teacher, special ed teacher, primary school teacher). Images for the 'pensioner' profession typically depict elderly individuals, which might be typical of the occupation but not necessarily sufficient to determine it from an image. We also observed that images generated of football player depict American football for White, Black and Latino races, while American soccer is depicted for this occupation for Indian, Asian and Middle Eastern races. We hypothesize that this is due to the fact that American soccer is often referred to as football in countries which are predominately populated by the latter three races, which influences how Stable Diffusion depicts this occupation for those races.

The depiction of some occupations were observed to contain intrinsic biases related to physical characteristics. For example, a majority of the images analyzed for the computer programmer and software developer professions depicted individuals wearing glasses. For umpires, security guards, opera singers, we observed that most of the images depicted individuals who were perceived to be overweight or obese. Images generated for farmers were perceived as having darker skin tones across all races, indicating the greater sun exposure for individuals with occupation. While not rep-

*a Middle Eastern female construction worker*

*a picture of an Asian female roofer.*

*a White male nurse*

*an image of a Black female chief executive officer*

Table 7. Examples of failure cases identified by manual error analysis

| Error Category | White | Black | Indian | Asian | M.E. | Latino |
|---|---|---|---|---|---|---|
| Failure to generate female subject | 3% | 4% | 1% | 2% | 12% | 5% |
| Failure to generate male subject | 4% | 2% | 1% | 1% | 0% | 2% |

Table 8. Error analysis with 100 samples for each race-gender combination

resentative of all possible biases, these observations indicate the presence of some intrinsic bias associated with certain

| gender | precision | recall | f1-score | support |
|--------|-----------|--------|----------|---------|
| male   | 0.99      | 0.95   | 0.97     | 549     |
| female | 0.96      | 0.99   | 0.97     | 589     |

Table 9. Confusion Matrix for male and female detection.

physical characteristics and occupations.

Similarly, we observed certain gender biases associated with occupations. For example, Stable Diffusion failed to generate male subjects across all races for the midwife and maid professions in the random subset of examples we analyzed. This could could be due to how these occupations are closely associated with the female gender, which makes it challenging for Stable Diffusion to generate male counterfactual images. However, no such bias was observed while generating images of policeman or handyman for female subjects across various races. We observed some shortcomings in the depictions of Middle Eastern and Latino female priests, possibly indicating a strong association between this occupation and the male gender for certain races.

Finally, for images which depicted multiple individuals, we occasionally observed a mix-up in the race associated with the occupation. For instance, images depicting hairdressers and barbers often associated the described race with the client rather than the hairdresser or barber. This is likely due to challenges that Stable Diffusion has with accurately binding attributes when multiple subjects are present in an image.

### 7.4. Generating counterfactuals for other subjects and social attributes

In addition to generating counterfactual sets depicting occupations, we also explored the viability of using personality traits as the subject in our captions. Specifically, we used the same list of 63 personality traits as in Naik and Nushi [38] to construct captions in a similar manner as described in Section 3.2, using the template $<p> <s> <a_1> <a_2>$. For example, given the personality trait *helpful*, we constructed captions for investigating race-gender bias such as *A helpful white man*, *A helpful black woman*, etc. However, we determined after a manual evaluation of generated images that Stable Diffusion struggled to depict these personality traits, which in turn degraded the quality of depictions of investigated social attributes. We therefore decided to omit images for these types of subjects from SocialCounterfactuals.

Beyond the race, gender, and physical characteristic social attributes, we also investigated the use of religion as a fourth type of social attribute for our counterfactual sets. This social attribute set included the terms *Christian, Muslim, Hindu, Buddhist, Atheist*, and *Agnostic*. After manual evaluation of images generated using religion as a social attribute, we determined that several of the religion terms consistently

produced nearly identical images (e.g., *Christian, Atheist, Agnostic*). Among images that produced discernible differences, Stable Diffusion primarily used race to differentiate subjects for each religion. Therefore, we decided to exclude counterfactual sets involving the religion social attribute from SocialCounterfactuals, and leave further investigation into intersectional biases involving religion to future work.

## 8. Dataset Generation & Experiment Details

### 8.1. Details of compute infrastructure used

The counterfactual image-text data was created using a large AI cluster equipped with Intel 3rd Generation Intel® Xeon® processors and Intel® Gaudi-2® AI accelerators. Up to 256 Intel Gaudi-2® AI accelerators were used to generate our SocialCounterfactuals dataset.

### 8.2. Details of training experiments

We fine-tune CLIP with learning rates of 7e-6, 5e-6 and 9e-6 for (Race, Gender), (Physical Char., Gender) and (Physical Char., Race) respectively for 1 epoch and a batch size of 32. For the FLAVA and ALIP, we follow the same fine-tuning setting as CLIP.

### 8.3. Details of dataset filtering

Table 11 provides details of the quantity of counterfactual sets remaining after each stage of filtering with our methodology. The total number of original counterfactual sets was 312,000, which corresponds to 5,408,000 generated images. After applying the CLIP similarity filter, the dataset was reduced to 21,359 counterfactual sets. This stage of filtering removes a significant portion of our generated dataset, but helps ensure that only images with the highest quality are retained. Once the NSFW filter was applied, the number of counterfactual sets decreased to 21,116, representing a 1% reduction. Finally, after the CLIP attribute detectability filter was applied, the dataset was reduced to 13,824 counterfactual sets, which is the final size of SocialCounterfactuals reported in Table 1. Below we provide additional details of each of the filtering stages.

**NSFW Filtering.** A manual analysis of the generated images detected NSFW samples. In order to detect and discard these images, we used an off-the-shelf fine-tuned vision transformer (google/vit-base-patch16-224-in21k) trained for NSFW image classification[9]. This model was fine-tuned used a proprietary dataset containing 80,000 images which was carefully curated to represent two classes: *nsfw* and *normal*. The reported evaluation accuracy of the fine-tuned model is 98.03.

---

[9]https://huggingface.co/Falconsai/nsfw_image_detection

After applying the NSFW classification model to our dataset, approximately 0.9% of the counterfactual sets were discarded for the (Race, Gender) segment, 1.4% for (Physical Char., Gender) segment, and 2.7% for (Physical Char., Race) segment. See Table 11 for details.

**CLIP Attribute Detectability Filtering.** To further ensure the quality of our dataset, we additionally filter counterfactual sets based on the ability of CLIP to detect and discern the targeted social attributes in each image. Intuitively, a counterfactual set is filtered out with respect to an attribute type if the number of its images that lack of detectability of that attribute type is less than a learnt threshold. Such a threshold is learnt based on the manual labels of 100 randomly sampled respective counterfactual sets. Then a counterfactual set is eventually filtered out if it is filtered out with respect to some of its attributes.

In more detail, we employed a two-phase approach for CLIP attribute detectability filtering. Without the loss of generality, considering an arbitrary intersectional bias $(X, Y)$ and a target attribute type X. Initially, we randomly sampled 100 counterfactual sets from $(X, Y)$-segment of our dataset. In the first phase, we manual labeled whether each counterfactual set is filtered out based on how many of its images possess their attribute type $X$ discernibly.

In the second phase, we develop a learnable threshold-based heuristic to filter out a counterfactual set if the number of its images that have their targeted attribute $X$ discernible by CLIP-based scores is less than the respective threshold. In particular, we use the names of all the attributes of the targeted attribute type to construct the set of potential text pairings as { a/an $x$ person} where $x$ is an attribute belonging to the attribute type $X$ and predict the most probable (image, text) pair according to CLIP-based image-text similarity scores. For instance, if $X$ is gender, then the potential text pairing set is {a female person, a male person}. An image is said to have their targeted attribute $x \in X$ discernible by CLIP-based score if the pair of the image and the text 'a/an $x$ person' is the most probable one. The learnable threshold were heuristically derived to obtain high correspondence between automatic filtering with CLIP-based scores and those filtered by the manual human annotation in the first phase. Table 10 lists learnable thresholds for each attribute pair and its respective attributes.

We acknowledge that in spite our best efforts, it is possible this manual analysis could propagate the annotator's bias associated with the various social attributes we investigate. While we also acknowledge that automatic filtering may introduce additional bias from CLIP, our error analysis in Section 7.2 shows that this filtering increases the quality of our dataset (90.8% → 97.5%). Additionally, our quantitative results (Table 2) show significant debiasing of other models (e.g., FLAVA) after training on our dataset, even when

| Attribute Pair | Threshold | | |
|---|---|---|---|
| | Gender | Race | Physical Char. |
| (Race, Gender) | 12 | 9 | N/A |
| (Physical Char., Gender) | 10 | N/A | 5 |
| (Physical Char., Race) | N/A | 13 | 8 |

Table 10. CLIP attribute detectability filtering thresholds for each attribute type in each attribute pair. N/A means not applicable when an attribute type is not a part of the attribute pair.

measured using other non-synthetic datasets (Table 3). This suggests that any additional bias introduced by our filtering method is less significant than the overall debiasing effect produced by training on our dataset.

Despite the impressive performance of text-to-image diffusion models, our error analysis (Section 7.2) shows they cannot be relied upon in an automated synthetic image generation pipeline without the use of filtering. Manual filtering by humans is not practical when generating a dataset at our scale (over 5.4 million images before filtering). Furthermore, counterfactual images depicting various combinations of intersectional social biases do not exist in natural image datasets. Therefore, we believe the use of automated filtering is necessary to construct a dataset which is useful for investigating intersectional social biases at scale.

### 8.4. Details of Caption Construction

We use a set of 260 occupations in this work, which was collected by combining the occupation lists proposed by Nadeem et al. [37], Chuang et al. [11], Naik and Nushi [38], and Harrison et al. [24]. After applying our three-stage filtering methodology, only 158 occupations remain in our dataset, which are provided in Table 18. To study bias associated with physical characteristics, we used the following keywords for positive and negative body stereotypes provided in Mei et al. [35]: skinny, obese, young, old, tattooed.

Examples of captions constructed using various prefixes, subjects, and bias attributes are provided in Table 19. We provide details of the total number of captions and images generated for each subject and attribute pairs in Table 11. The total number of counterfactual sets is determined by the product of the number of prefixes used to construct captions (4), the cardinality of the subject set (i.e., the number of occupations), and 100 (the number of over-generations per set). The number of images per set is determined by the product of the cardinalities of the attribute sets. The total number of generated images is the product of the number of counterfactual sets and the number of images per set.

| Attribute Types | No. Sets | Images Per Set | No. Images | Sets After CLIP Sim. Filter | % Filtered Out | Sets After NSFW Filter | % Filtered Out | Sets After CLIP Attrib. Filter | % Filtered Out | Final No. Images |
|---|---|---|---|---|---|---|---|---|---|---|
| (Race, Gender) | 104,000 | 12 | 1,248,000 | 13,147 | 87% | 13,035 | 0.9% | 7,936 | 39% | 95,232 |
| (Physical Char., Gender) | 104,000 | 10 | 1,040,000 | 7,254 | 93% | 7,149 | 1.4% | 5,052 | 29% | 50,520 |
| (Physical Char., Race) | 104,000 | 30 | 3,120,000 | 958 | 99% | 932 | 2.7% | 836 | 10% | 25,080 |
|  | 312,000 |  | 5,408,000 | 21,359 |  | 21,116 |  | 13,824 |  | 170,832 |

Table 11. Details of the number of counterfactual sets after applying different filters in each of the stages.

# 9. Additional Discussion of Limitations and Ethical Considerations

**Limitations** Despite our best efforts, the templates and methodologies we adopt may themselves contain some latent biases which could contribute to the implicit biases exhibited by VLMs. All statements pertaining to gender, race, and occupational attributes or associations should be interpreted only within the context of our experiments. Furthermore, all discussion of social attributes in this work are intended to be interpreted as *perceived*. We are aware that our approach only considers binary classification of genders and does not exhaustively encompass all races, physical characteristics, and occupations, which is due to data limitation rather than value judgements by the authors. The results we present cannot be assumed to generalize to social and demographic terms omitted in our analysis. The labels for the attributes we present in the paper are derived from prior work and were further limited to those which stable diffusion could depict. Our goal is to provide text labels that produce perceived physical differences, but these are not labels we aim to impose on any groups or sub-groups. Similar to Smith et al. [49], we recognize there are trade-offs in creating lists of socials attributes. While these lists may not be entirely inclusive, we leverage them for their benefit in identifying and mitigating bias. Our study was conducted in English, which limits the generalizability of our findings to other languages.

**Ethical Considerations.** With the findings we present in this paper, we aim to increase the understanding of bias in VLMs and probe mitigation strategies. We acknowledge that our work does not encompass all possible social attributes and that our selected categories for gender, race, physical characteristics, and occupations may harbor stereotypes that cannot be assumed to represent their entire groups. Similar to Hall et al. [23], we recognize that we may miss intersectional characteristics that constitute a well-accepted image of a person in a specific occupation or belonging to a race. Our aim is that the techniques presented in this work can help reduce various social disparities in VLMs and can be further extended to include more genders, races, occupations and other social characteristics. Continuing these efforts will increase confidence in the ability of VLMs to exhibit fairness with respect to differing social attributes. We understand

| Model | Mean | Std | Min | 25% | 50% | 75% | Max |
|---|---|---|---|---|---|---|---|
| ALIP | 0.58 | 0.23 | 0.28 | 0.41 | 0.51 | 0.67 | 1.31 |
| CLIP | 0.34 | 0.09 | 0.24 | 0.29 | 0.32 | 0.37 | 0.76 |
| FLAVA | 0.34 | 0.08 | 0.23 | 0.28 | 0.32 | 0.39 | 0.57 |
| LACLIP | 0.73 | 0.26 | 0.29 | 0.50 | 0.70 | 0.90 | 1.59 |
| OPENCLIP | 0.77 | 0.30 | 0.31 | 0.51 | 0.75 | 0.93 | 2.04 |
| SLIP | 0.71 | 0.28 | 0.26 | 0.47 | 0.65 | 0.89 | 1.41 |

Table 12. Distribution of NDKL scores by model, measured across occupations using the **race-gender** dataset segment.

| Model | Mean | Std | Min | 25% | 50% | 75% | Max |
|---|---|---|---|---|---|---|---|
| ALIP | 0.52 | 0.25 | 0.21 | 0.33 | 0.43 | 0.73 | 1.05 |
| CLIP | 0.47 | 0.29 | 0.14 | 0.19 | 0.36 | 0.67 | 1.03 |
| FLAVA | 0.45 | 0.27 | 0.13 | 0.20 | 0.38 | 0.66 | 0.92 |
| LACLIP | 0.56 | 0.25 | 0.24 | 0.36 | 0.47 | 0.77 | 1.05 |
| OPENCLIP | 0.54 | 0.25 | 0.21 | 0.32 | 0.46 | 0.71 | 1.05 |
| SLIP | 0.56 | 0.26 | 0.23 | 0.34 | 0.45 | 0.76 | 1.19 |

Table 13. Distribution of NDKL scores by model, measured across occupations using the **race-physical characteristics** dataset segment.

that the use of a bias reduction strategy without deep understanding of various nuances does not guarantee a foolproof solution in bias elimination, and still may result in VLMs that cause harm and stigmatize certain subsets of individuals. Therefore, debiasing efforts should be further developed prior to wide-spread adoption in sensitive applications.

## 10. Additional Analysis and Results

### 10.1. NDKL and Bias@K results

In addition to MaxSkew, we also calculated the Bias@K and Normalized Discounted Kullback-Leibler Divergence (NDKL) metrics for our bias probing & mitigation experiments. See Geyik et al. [20] for a detailed description of NDKL and Wang et al. [52] for details of Bias@K.

We estimate NDKL by summing over ranked lists of size $\{1, 2, .., K^2\}$, where $K = |A_1| \times |A_2|$. Tables 12, 13, and 14 provide the distribution of NDKL scores for each model, measured across occupations for the three segments
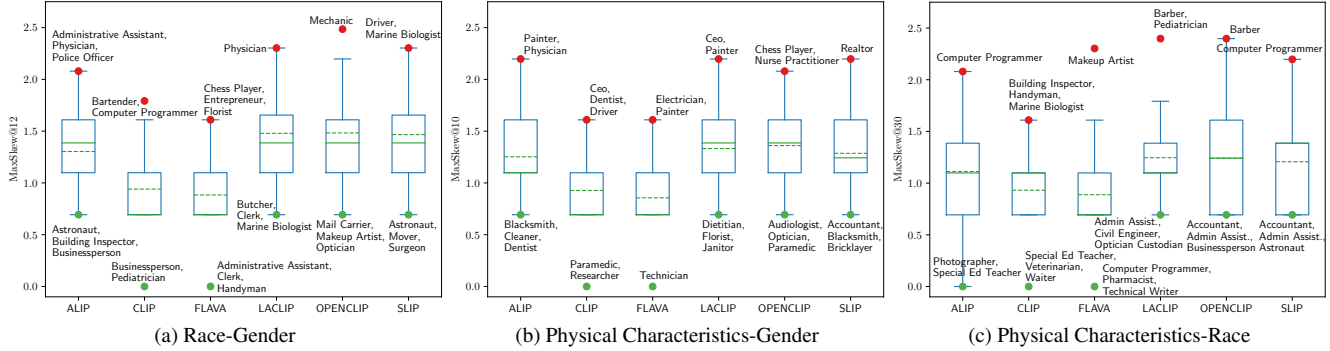
Figure 6. Distribution of MaxSkew@$K$ measured across occupations for (a) Race-Gender, (b) Physical Characteristics-Gender, and (c) Physical Characteristics-Race intersectional biases after NSFW and Attribute Detectability Filtering. Max (min) values are plotted as red (green) circles with corresponding occupation names

| Model | Mean | Std | Min | 25% | 50% | 75% | Max |
|---|---|---|---|---|---|---|---|
| ALIP | 0.63 | 0.28 | 0.28 | 0.43 | 0.54 | 0.76 | 1.83 |
| CLIP | 0.39 | 0.08 | 0.27 | 0.34 | 0.38 | 0.43 | 0.65 |
| FLAVA | 0.36 | 0.07 | 0.27 | 0.31 | 0.34 | 0.39 | 0.55 |
| LACLIP | 0.72 | 0.29 | 0.29 | 0.49 | 0.66 | 0.91 | 1.68 |
| OPENCLIP | 0.74 | 0.26 | 0.32 | 0.52 | 0.68 | 0.91 | 1.55 |
| SLIP | 0.66 | 0.24 | 0.31 | 0.49 | 0.60 | 0.79 | 1.28 |

Table 14. Distribution of NDKL scores by model, measured across occupations using the **gender-physical characteristics** dataset segment.

| Dataset Segment | ALIP | CLIP | FLAVA | LACLIP | OPENCLIP | SLIP |
|---|---|---|---|---|---|---|
| Physical-Gender | 0.06 | -0.18 | **0.02** | 0.51 | 0.50 | 0.40 |
| Race-Gender | **0.01** | -0.12 | -0.02 | 0.47 | 0.50 | 0.41 |

Table 15. Bias@K by model, calculated across occupations for dataset segments which include gender attributes.

of SocialCounterfactuals. Consistent with our analysis of MaxSkew@K (Section 4.2), CLIP and FLAVA generally have the lowest NDKL scores. Notably, all models exhibit values well above the ideal case of 0 for this metric.

Table 15 provides Bias@K for our two dataset segments which include gender attributes, where we set $K = |A_1| \times |A_2|$ as in our MaxSkew@K and NDKL evaluations. Bias@K is a measure of marginal gender bias; a value of 0 indicates that both genders are represented equally in retrieval results. Positive values indicate that males are represented more than females, while negative values indicate that females are represented more than males. ALIP and FLAVA exhibit the least amount of gender bias overall. CLIP demonstrates some bias towards female images, while LACLIP, OpenCLIP, and SLIP all exhibit a strong bias for male images.

Tables 16 and 17 provide NDKL and Bias@K results (re-

spectively) for models which were evaluated in our debiasing experiments (Section 5). Because Bias@K only measures gender bias, we only calculate it for dataset segments which include gender attributes.

## 10.2. Preliminary evaluations using counterfactuals with three intersectional attributes

Our dataset generation approach can incorporate greater combinations of attributes by introducing additional attribute sets to our caption templates (Section 3.2). To investigate this, we produced 16k images spanning 266 counterfactual image sets with three intersectional attributes in each caption (physical characteristics, gender, and race). The CLIP MaxSkew@60 score for three-attribute intersectional bias on this set of images is 0.693, which reveals even greater bias than when measuring various pairs of two-attribute intersectional bias (which have a MaxSkew@60 of 0.567-0.59 on the same images).

We believe this preliminary result points to a promising direction for future studies, which could apply our approach to investigate bias across a greater number of intersectional attributes. To the best of our knowledge, our work is the first to address the problems of probing and mitigating intersectional bias in vision-language models. While we focused on pairs of attributes in our paper, we believe this nonetheless represents a significant contribution considering the lack of prior work addressing intersectional social biases.

## 10.3. FID and IS evaluation of SocialCounterfactuals and other datasets

We computed Inception Score (IS) and FID relative to ImageNet and LAION-2B for our dataset, VisoGender, and PATA. As shown in the table below, our dataset outperforms both VisoGender and PATA (which consist of real-world images for bias evaluation) across all metrics.

| Intersectional Bias | CLIP [41] | | ALIP [59] | | FLAVA [46] | |
|---|---|---|---|---|---|---|
| | Pre-trained | Debiased | Pre-trained | Debiased | Pre-trained | Debiased |
| (Race, Gender) | 0.37 | 0.33 | 0.71 | 0.36 | 0.376 | 0.365 |
| (Physical Char., Gender) | 0.41 | 0.34 | 0.77 | 0.36 | 0.38 | 0.40 |
| (Physical Char., Race) | 0.20 | 0.17 | 0.30 | 0.20 | 0.204 | 0.196 |

Table 16. Mean of NDKL for pre-trained and debiased variants of CLIP, ALIP, and FLAVA, estimated by withholding counterfactual sets for 20% of the occupations in our dataset.

| Intersectional Bias | CLIP [41] | | ALIP [59] | | FLAVA [46] | |
|---|---|---|---|---|---|---|
| | Pre-trained | Debiased | Pre-trained | Debiased | Pre-trained | Debiased |
| (Race, Gender) | -0.19 | 0.13 | 0.15 | 0.20 | -0.05 | 0.10 |
| (Physical Char., Gender) | -0.22 | 0.06 | 0.09 | 0.03 | -0.09 | 0.02 |

Table 17. Bias@K for pre-trained and debiased variants of CLIP, ALIP, and FLAVA, estimated by withholding counterfactual sets for 20% of the occupations in our dataset.

| | FID (ImageNet) ↓ | FID (LAION-2B) ↓ | IS ↑ |
|---|---|---|---|
| VisoGender | 130.44 | 113.45 | 6.2 |
| PATA | 116.29 | 98.34 | 11.3 |
| SocialCounterfactuals | **106.83** | **89.91** | **12.07** |

## 10.4. Additional details and results for bias probing experiments

Our bias probing analysis presented in Section 4.2 utilized our SocialCounterfactuals dataset before NSFW and Attribute Detectability Filtering was applied. For completeness, in this section we provide the same analysis on the SocialCounterfactuals dataset after NSFW and Attribute Detectability Filtering.

Figure 6 provides boxplots of the intersectional bias MaxSkew@$K$ distribution for each VLM, measured across occupations separately using the three segments of our dataset after NSFW and Attribute Detectability Filtering. Overall we find that these distributions are largely similar to those described previously in Figure 3. Similarly, the marginal gender skewness depicted in Figure 7 and retrieval proportions for the 'Doctor' profession illustrated in Figure 8 (post NSFW and Attribute Detectability Filtering) reflect largely the same trends as those discussed previously in Section 4.2 for Figures 4 and 5. We also provide retrieval proportions for other occupations in Figures 9 to 13.
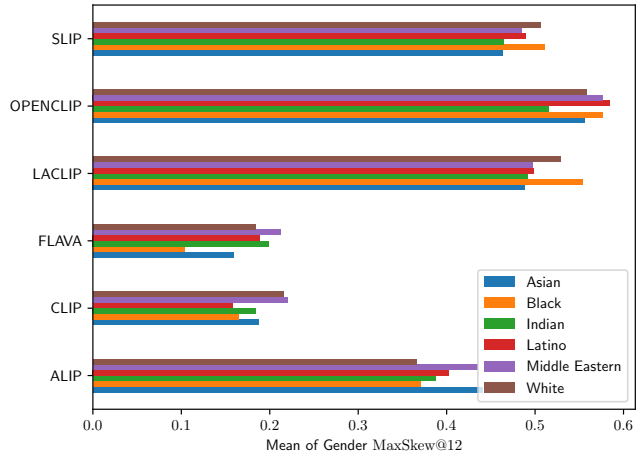


Figure 7. Mean of (marginal) gender MaxSkew@$K$ measured across occupations for different races after NSFW and Attribute Detectability Filtering.
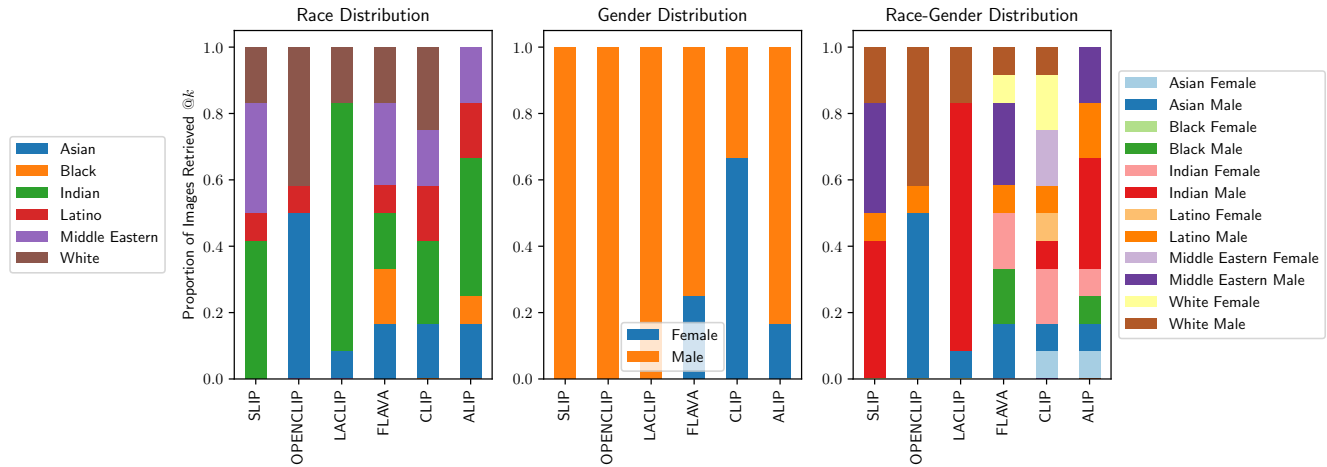
Figure 8. Proportion of images retrieved @$k = 12$ after NSFW and Attribute Detectability Filtering using neutral prompts for the **Doctor** occupation, broken down by race & gender attributes.
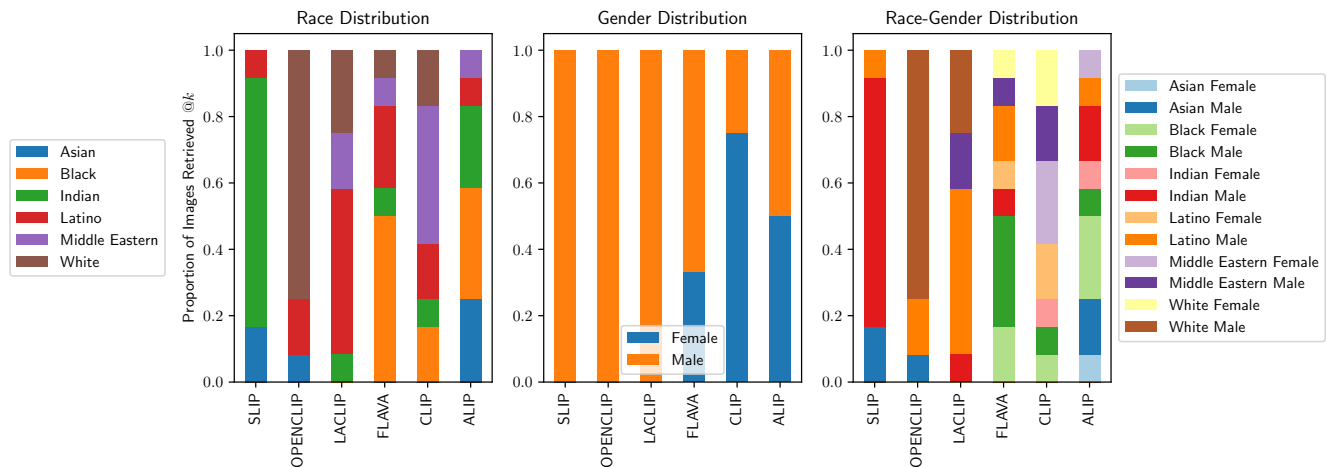


Figure 9. Proportion of images retrieved @$k = 12$ after NSFW and Attribute Detectability Filtering using neutral prompts for the **Software Developer** occupation, broken down by race & gender attributes.
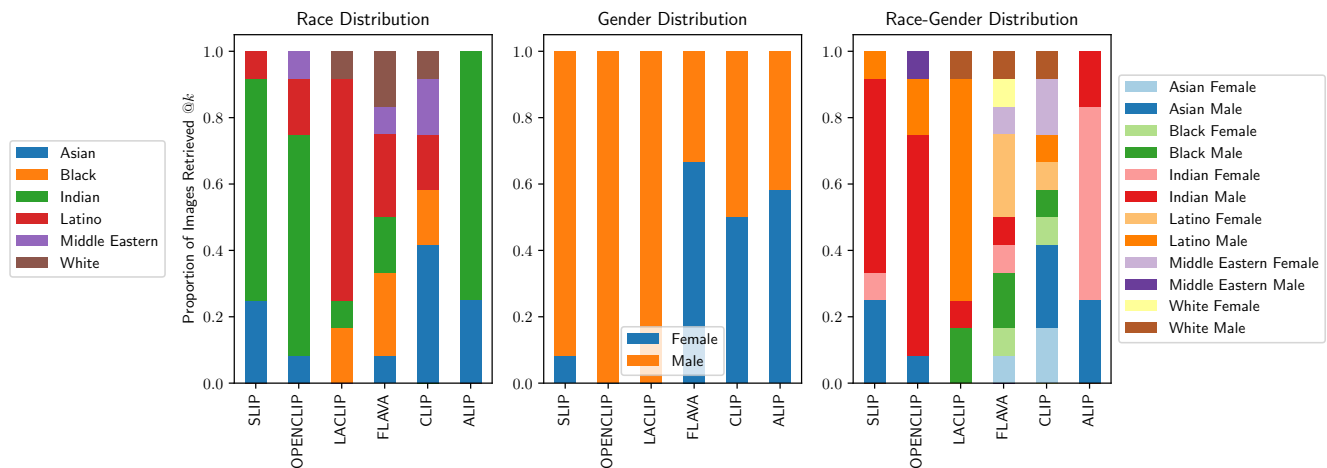


Figure 10. Proportion of images retrieved @$k = 12$ after NSFW and Attribute Detectability Filtering using neutral prompts for the **Construction Worker** occupation, broken down by race & gender attributes.
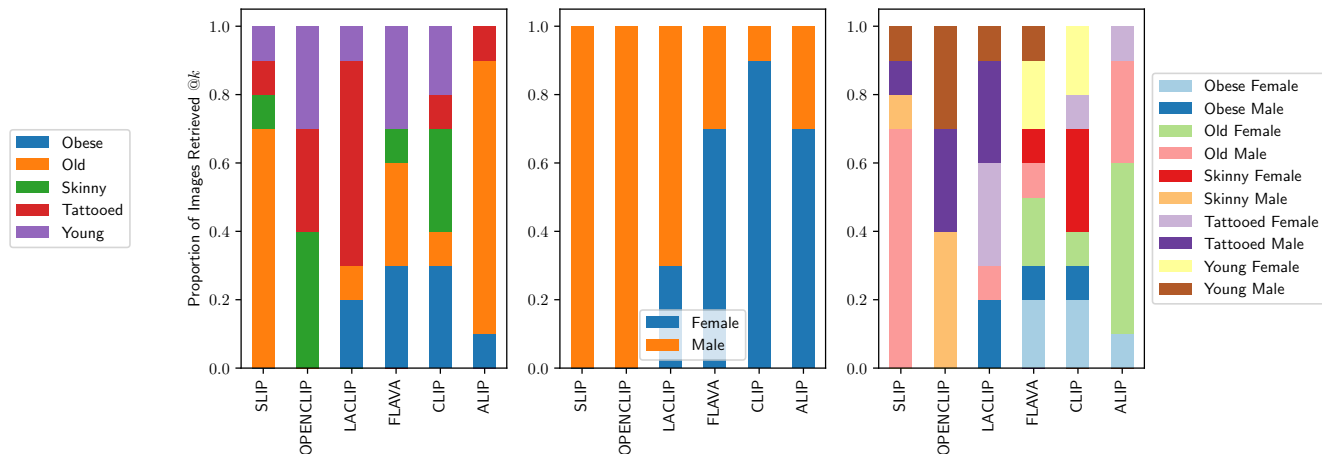
Figure 11. Proportion of images retrieved @$k = 10$ after NSFW and Attribute Detectability Filtering using neutral prompts for the **Entrepreneur** occupation, broken down by gender & physical characteristics.
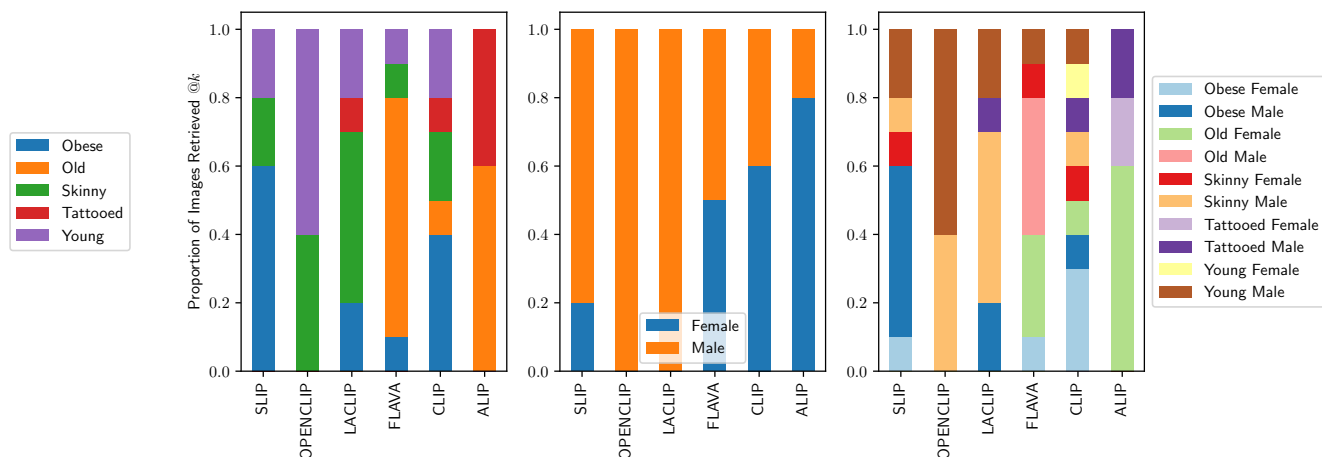


Figure 12. Proportion of images retrieved @$k = 10$ after NSFW and Attribute Detectability Filtering using neutral prompts for the **Technical Writer** occupation, broken down by gender & physical characteristics.
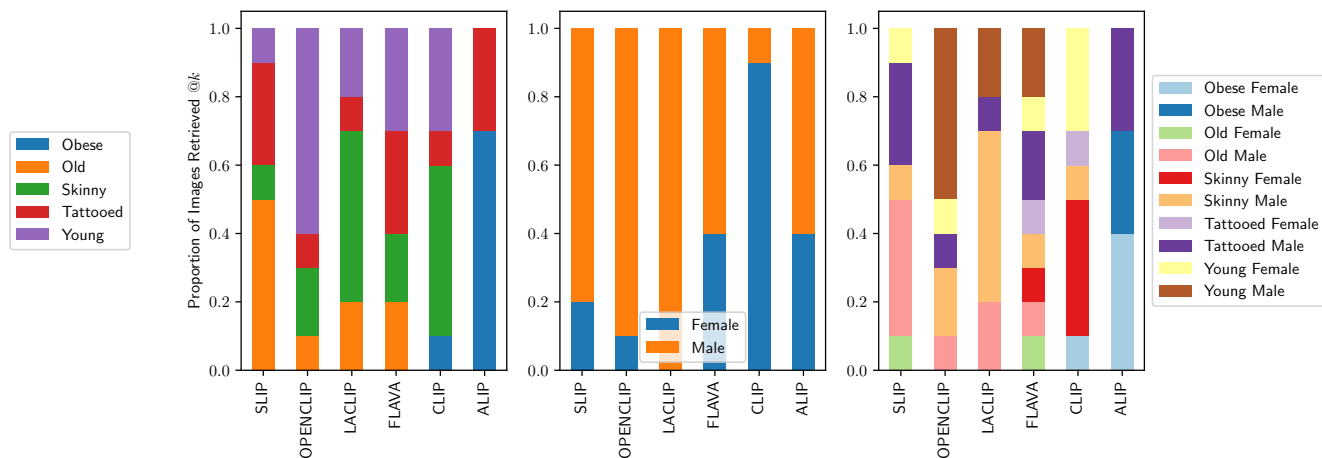


Figure 13. Proportion of images retrieved @$k = 10$ after NSFW and Attribute Detectability Filtering using neutral prompts for the **Salesperson** occupation, broken down by gender & physical characteristics.

**Occupations**

academic, accountant, administrative assistant, analyst, architect, army, artist, assistant, astronaut, athlete, attorney, audiologist, auditor, author, baker, banker, barber, bartender, biologist, blacksmith, boxer, bricklayer, broker, building inspector, bus driver, businessperson, butcher, carpenter, cashier, ceo, chef, chemist, chess player, chief, chief executive officer, childcare worker, civil engineer, civil servant, cleaner, clerk, coach, comedian, commander, composer, computer programmer, construction worker, consultant, cook, crane operator, customer service representative, dancer, delivery man, dentist, designer, detective, dietitian, dj, doctor, driver, economist, editor, electrician, engineer, entrepreneur, farmer, firefighter, florist, football player, graphic designer, guard, guitarist, hairdresser, handball player, handyman, housekeeper, janitor, judge, lab tech, laborer, lawyer, librarian, magician, mail carrier, makeup artist, manager, marine biologist, mathematician, mechanic, model, mover, musician, nurse, nurse practitioner, nutritionist, opera singer, optician, optician custodian, painter, paramedic, pastry chef, pediatrician, pensioner, pharmacist, photographer, physician, physicist, pianist, pilot, plumber, poet, police officer, policeman, pr person, priest, primary school teacher, prisoner, producer, professor, psychologist, real estate developer, real estate agent, realtor, receptionist, recruiter, reporter, researcher, roofer, sailor, salesperson, scientist, secretary, security guard, sheriff, software developer, soldier, special ed teacher, statistician, supervisor, surgeon, surveyor, swimmer, tailor, teacher, technical writer, technician, telemarketer, tennis player, therapist, tour guide, umpire, veterinarian, videographer, waiter, web developer, writer, zoologist

Table 18. Full list of occupations which remain in our dataset after filtering.

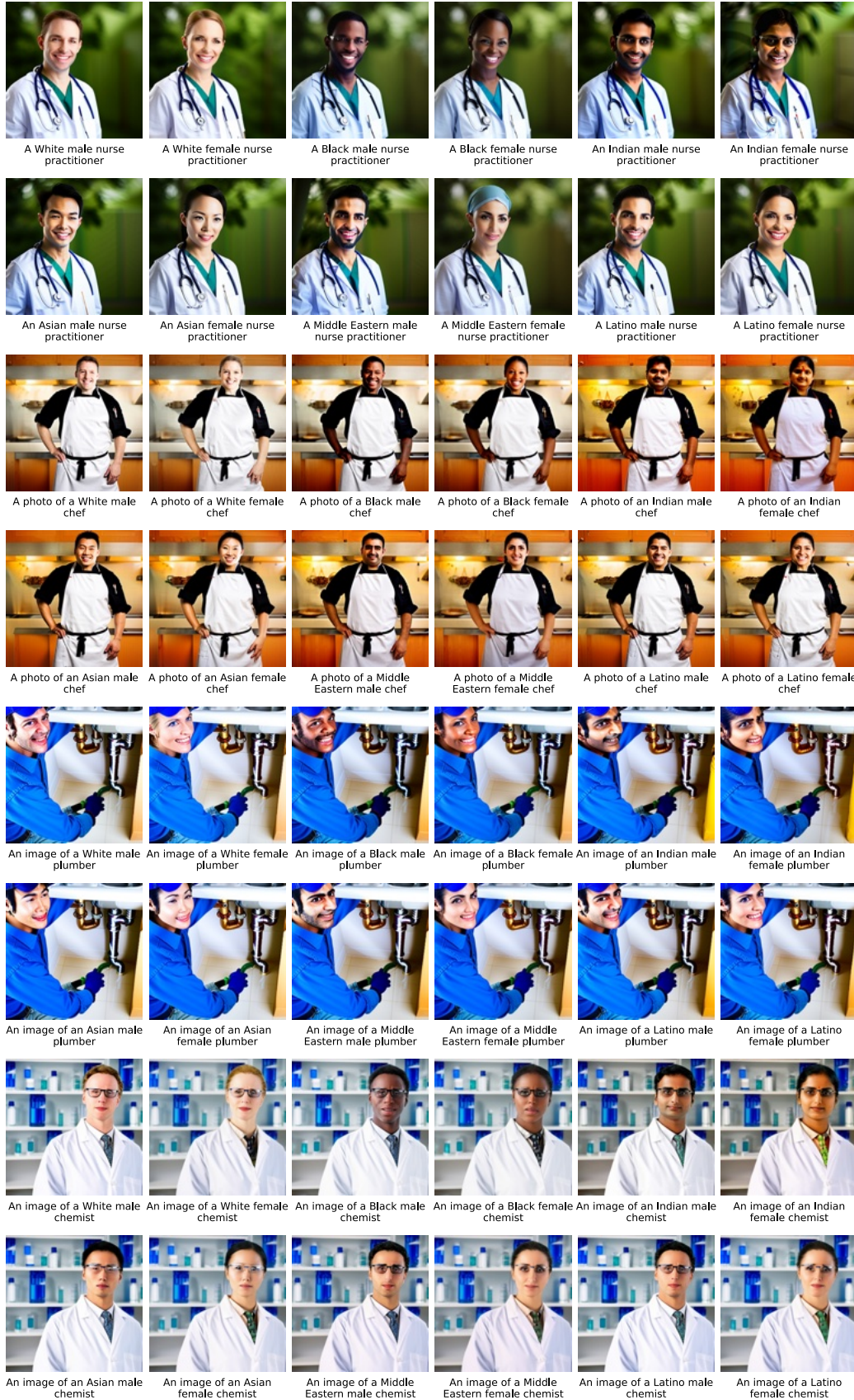| Template | Attribute-1 | Attribute-2 | Subject | Caption |
|---|---|---|---|---|
| A/An <race> <gender> <occupation> | Indian | female | nurse | An Indian female nurse |
| A photo of <race> <gender> <occupation> | White | male | electrician | A photo of a White male electrician |
| A picture of <physical char.> <race> <occupation> | young | Latino | web developer | A picture of a young Latino web developer |
| An image of <physical char.> <gender> <occupation> | tattooed | male | barber | An image of a tattooed male barber |

Table 19. Examples of captions constructed from templates

Figure 14. Additional examples of counterfactual sets produced by our approach for the (Race, Gender) attribute pair.

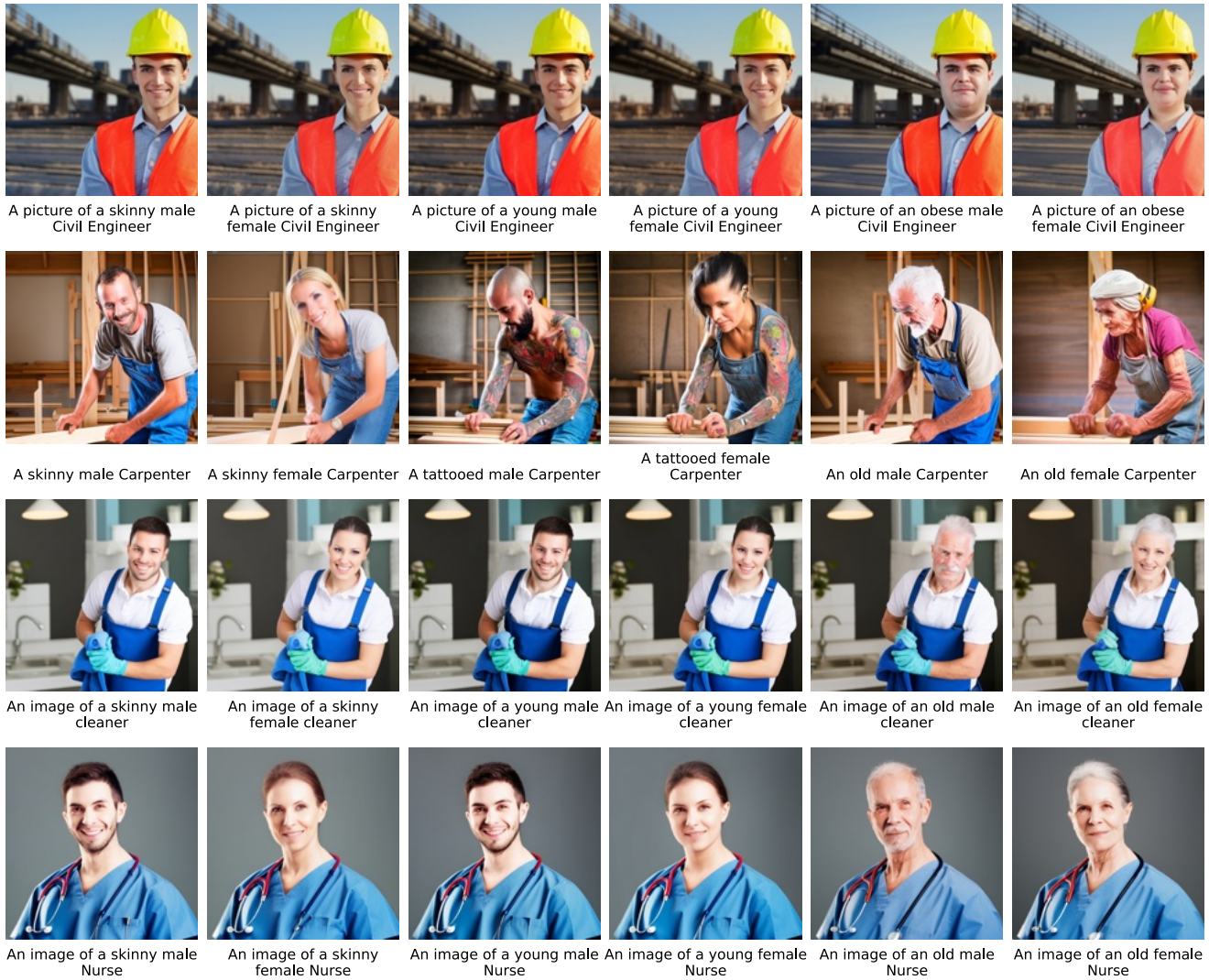| | | | | | |
|---|---|---|---|---|---|
| A picture of a skinny male Civil Engineer | A picture of a skinny female Civil Engineer | A picture of a young male Civil Engineer | A picture of a young female Civil Engineer | A picture of an obese male Civil Engineer | A picture of an obese female Civil Engineer |
| A skinny male Carpenter | A skinny female Carpenter | A tattooed male Carpenter | A tattooed female Carpenter | An old male Carpenter | An old female Carpenter |
| An image of a skinny male cleaner | An image of a skinny female cleaner | An image of a young male cleaner | An image of a young female cleaner | An image of an old male cleaner | An image of an old female cleaner |
| An image of a skinny male Nurse | An image of a skinny female Nurse | An image of a young male Nurse | An image of a young female Nurse | An image of an old male Nurse | An image of an old female Nurse |

Figure 15. Additional examples of counterfactual sets produced by our approach for the (Physical Characteristics, Gender) attribute pair.

Figure 16. Additional examples of counterfactual sets produced by our approach for the (Physical Characteristics, Race) attribute pair.