

A Stereotype Content Analysis on Color-related Social Bias in Large Vision Language Models

Anonymous ACL submission

Abstract

As large vision language models (LVLMs) rapidly advance, concerns about the risk of their learning or generating stereotypes are increasing. However, previous studies on LVLM’s stereotypes face two primary limitations: metrics that overlooked the importance of content words, and datasets that overlooked the effect of color. To address these limitations, this study introduces new evaluation metrics based on the Stereotype Content Model (SCM). The SCM-based metric is grounded in social psychology, enabling the detection of stereotypical content words in model responses along the dimensions of competence and warmth. We also propose BASIC, a benchmark for assessing gender, race and color stereotypes. Using SCM metrics and BASIC, we conduct a study with eight LVLMs to discover stereotypes. As a result, we found three findings. (1) The SCM-based evaluation is effective in capturing stereotypes. (2) LVLMs exhibit color stereotypes in the output along with gender and race ones. (3) Interaction between model architecture and parameter sizes seems to affect stereotypes. We release BASIC publicly on [anonymized for review].

1 Introduction

As large vision-language models (LVLMs) continue to advance, they are increasingly being utilized in perceiving visual images across diverse domains (Xu et al., 2024; Gu et al., 2024; Wang et al., 2025). However, when employing LVLMs, it is essential to acknowledge that models may be affected by stereotypes, potentially influencing the models’ perception. In humans, a stereotype refers to cognitive generalizations about the characteristics of individuals belonging to a specific social group (VandenBos, 2007). So, people often judge others unconsciously based on their visual appearance such as gender, race, and color (Cloutier et al., 2005; Eberhardt et al., 2004; Elliot et al., 2013).

Given that LVLMs are usually trained on human-generated data (Jia and Liang, 2017; Cheng et al., 2019), it is likely that they have learned human stereotypes. As a result, LVLMs sometimes exhibit harmful stereotypes, leading to bad consequences (Jiang et al., 2024; Howard et al., 2025; Fraser and Kiritchenko, 2024). Therefore, the need to identify such stereotypes is becoming important. Researchers recently designed various benchmarks to quantify stereotypes in LVLMs (Adewumi et al., 2024; Wang et al., 2024b; Park et al., 2025).

Despite the success of uncovering stereotypes, existing studies have less considered two factors possibly correlated with stereotypes: content words and color. First, regarding content words, we suspect that LVLMs may use positive sentiment words when generating stereotypical content. Many previous studies have relied on sentiment analysis to assess stereotypes (Wang et al., 2024a; Raj et al., 2024; Gerych et al., 2024). While this approach has contributed, sentiment-based evaluation cannot fully mirror how humans evaluate such stereotypes. Stereotypes are complex constructs compared to perceived sentiment, and thus cannot be thoroughly examined by sentiment polarity alone (Greenwald and Banaji, 1995; Amodio, 2014; Howard et al., 2025; Siy and Cheryan, 2016). Hence, to evaluate stereotypes more precisely, we adopt Stereotype Content Model (SCM), a socio-psychological framework that captures such implicit stereotypes in text and complements sentiment analysis (Fiske et al., 2018). With SCM, we examine how LVLMs exhibit stereotypical content in their responses.

Second, regarding color, LVLMs may exhibit stereotypes because of different colors in the image. Previously, researchers in vision-language model reported that model responds inconsistently due to image differences, even when images contain similar scenes (Wu et al., 2025; Kantipudi et al., 2020; Liang et al., 2025). So, studies on stereotypes have employed image pairs to control such

083 differences. Previous attempts built controlled im-
 084 age pairs, however, they overlooked the potential
 085 impact of color tone. Color tone is critical because
 086 humans are influenced by tonal differences (Elliot
 087 and Maier, 2013; Kaya and Epps, 2004; Wexner,
 088 1954) and LVLMs might learn such tendency from
 089 human-generated data. Hence, we suspect that
 090 LVLm stereotypes could be affected by color tones.
 091 To justify such suspicion, we evaluate LVLMs with
 092 our new color-controlled benchmark.

093 Thus, we propose two metrics inspired by SCM
 094 and a new Benchmark for Assessing Stereotypes
 095 with Image Colors, BASIC. Additionally, we at-
 096 tempt to examine which part of LVLMs affect
 097 stereotypes, by comparing different models. We
 098 believe that a comparison between model architec-
 099 tures can hint at components relevant to stereotypes.
 100 So, this paper has the following contributions:

1. Based on SCM, we derive automatic metrics from human perception on stereotypes.
2. We suggest a benchmark that can separate the effect of colors from other stereotypes.
3. Comparing eight LVLMs, we found the effect of model and its size on stereotype is limited.

2 Related Work

108 We review prior study on evaluating stereotypes in
 109 terms of three aspects: (1) sentiment-based metrics,
 110 (2) designed benchmarks, and (3) tested LVLMs.

111 First, regarding content words, prior studies have
 112 mainly assessed model responses based on senti-
 113 ment polarity (Inoshita and Zhou, 2024; Hu et al.,
 114 2025; Gerych et al., 2024), which limits their abil-
 115 ity to capture the subtle meanings embedded in
 116 text. In the case of LVLMs, responses often ex-
 117 hibit positive sentiments because they were rein-
 118 forced by human feedback (Li et al., 2024), mak-
 119 ing it difficult to detect stereotypes solely through
 120 sentiment-based metrics. Therefore, more sophis-
 121 ticated measures that directly target stereotypes
 122 are needed. In response to this limitation, recent
 123 work has begun incorporating social-psychological
 124 evaluation frameworks to analyze stereotypes. For
 125 instance, (Howard et al., 2025) was influenced by
 126 the Stereotype Content Model (SCM). The study
 127 used the dictionary from (Nicolas et al., 2021) to
 128 measure the frequency of SCM-related words in im-
 129 age captions. However, this approach is limited by
 130 its reliance on word frequency, without accounting
 131 for contextual or derived meanings.



Figure 1: Examples from our BASIC dataset. Top row shows chess players with different colors (white, blue, red), retaining race and gender. Bottom row shows DJs with different race or gender, retaining white color.

132 Second, regarding benchmarks, prior studies
 133 have emphasized constructing paired images to
 134 assess stereotypes of LVLMs. Traditional stud-
 135 ies constructed diverse images based on common
 136 themes such as gender or race, and compared av-
 137 erage evaluation scores across these thematic sets
 138 (Mandal et al., 2023; Janghorbani and De Melo,
 139 2023; Qiu et al., 2023). However, such methods of-
 140 ten suffer from insufficient control over confound-
 141 ing variables. For example, if gender, race, and
 142 background change simultaneously, any observed
 143 stereotypes cannot be solely attributed to gender
 144 when examining gender stereotypes. So, recent
 145 works have proposed that changing a targeted at-
 146 tribute while retaining the rest of the image (Kim
 147 et al., 2025; Howard et al., 2025, 2024; Fraser and
 148 Kiritchenko, 2024). However, unintended color
 149 differences still appear in the generated images.
 150 Considering the visual differences can influence
 151 the model’s internal representations (De and Ped-
 152 ersen, 2021; Kantipudi et al., 2020; Liang et al.,
 153 2025), thereby potentially affecting the manifesta-
 154 tion of stereotypes. Thus, color differences should
 155 be considered in stereotype evaluation frameworks.

156 Lastly, regarding model architecture, researchers
 157 recently investigated whether the model architec-
 158 ture affects stereotypes. Though previous studies
 159 have identified stereotypes within a limited set of
 160 models (Raj et al., 2024; Sathe et al., 2024), they
 161 did not clearly mention which architectural choice
 162 cause stereotypes. Recently, researchers have be-
 163 gun to identify causal relationships (Wu et al., 2024;
 164 Wang et al., 2024a). Despite their success, it is yet
 165 questionable whether the findings could be gener-
 166 alized to large VLMs with more than 50 billion

parameters. So, we systematically incorporate a broader range of model sizes.

3 BASIC Benchmark

In this study, we aim to separately examine the effect of color in the analysis of stereotypes produced by LVLMs. There are various visual attributes in person images that can influence stereotype formation: gender, race, and color. Since these factors appear together in a single image, it is essential to separate the impact of each attribute on the model’s responses. To examine the potential impact of such differences, we construct a paired image dataset as shown in Figure 1.

We created the dataset using three primary colors: white, red, and blue. While a broader range of colors could have been included, we selected these particular ones based on cognitive factors related to human perception. Specifically, prior research has shown that humans tend to associate red, blue, and white with particular emotions, cognitive judgments, or stereotypes (Elliot and Maier, 2013; Kaya and Epps, 2004; Wexner, 1954). Therefore, we adapted the procedure proposed by (Howard et al., 2024) to construct a dataset with controlled color tones and designed a six-step process¹.

For the experimental analysis of model stereotypes, gender and race categorizations were adopted from previous studies (Karkkainen and Joo, 2021; Howard et al., 2024). Specifically, gender was determined based on biological sex, while race was classified into six predefined groups: African-American, Asian, Caucasian, Indian, Latino, and Middle Eastern. This selection was made solely for experimental purposes and is not intended to exclude any group.

Step 1: Constructing Occupation List. To create paired images, we started by choosing a common occupation for the individuals depicted in each image. Following the approach of Howard et al. (2024), we collected a list of occupations from studies (Nadeem et al., 2021; Chuang et al., 2023; Naik and Nushi, 2023). After removing duplicates, we obtained a list of 181 occupations. All subsequent steps are performed for each of 181 occupations.

Step 2: Creating action. We considered various actions that could be associated with each occupation. Since a single image cannot cover all possible

actions of a specific occupation, multiple images describing different actions are required for each occupation. To achieve this, we adapted a prompt from (Sathe et al., 2024) and used GPT-4 to generate five distinct scenarios for each occupation.

Step 3: Generating Seed Images. For each of the five actions associated with a given occupation, we generated a seed image as a baseline for transformation. To ensure stability in subsequent image manipulation steps, the seed image was generated using an Asian male wearing white clothing on a gray background. This step was performed using Stable Diffusion XL, which generated high-resolution images (1024 × 1024 pixels) based on scenario corresponding to each occupation (Podell et al., 2023). We employed negative prompts (Ban et al., 2024) to enhance image quality and ensure alignment with the intended depiction. A negative prompt is an input method that specifies unwanted attributes, guiding the model to avoid including them in the generated images.

Step 4: Controlling Gender and Race. Based on a seed image, we generate paired images that reflect various combinations of gender and race. In this stage, we utilize the Image-to-Image model of Stable Diffusion XL to modify the seed image according to the specified race and gender attributes. We retain the contextual information conveyed by the scenario while altering only the race and gender of the original image. To ensure that the pose, position, and background remain consistent with the seed image, we adopt negative prompts.

Step 5: Controlling Color. Based on the images generated in Step 4, we further modify the color tone of each image to either blue or red, while maintaining the same race and gender. To ensure that the intended color is clearly reflected in the output, we use negative prompts. Through this step, we obtain a total of 180 images per scenario (= 5 actions × 2 genders × 6 races × 3 colors).

Step 6: Filtering. To ensure the quality of the final dataset, we conducted a filtering process based on three criteria. First, we evaluated the semantic similarity between the 32,580 generated images and their corresponding text prompts. We normalized the embeddings of both image and text using CLIP, then computed the cosine similarity between them. Only images with a similarity score of 0.2 or higher were retained in the dataset. Second, we included only those images that passed the NSFW

¹The generation code is available at [anonymized for the review]. Detailed prompts and parameter settings for each step are provided in Appendix A.

	<i>Red</i>	<i>Green</i>	<i>Blue</i>
Red vs. White	0.12	-3.06	-2.43
Blue vs. White	-5.37	-4.87	-0.52

Table 1: RGB value changes (%) in red/blue images compared to white seed image

filtering, thereby excluding any potentially unsafe content that may have been generated. Lastly, we retained only occupations for which all 180 image variations passed both filtering conditions. If any single image was filtered out, we excluded the entire occupation to ensure completeness of pairs.

Through the above procedures, we generated a total of 18,360 images across 102 occupations. After completing all image generations, we assessed whether the color tones were clearly distinguishable by two methods. First, we computed the RGB scale differences between the red or blue images and their corresponding white counterparts. As shown in Table 1, the constructed dataset actively suppresses non-target colors to emphasize the intended color tone. Second, we manually validated 10% random samples to ensure the quality of dataset. We detail the result in Appendix B.

4 Experimental Setup

Using the dataset, we experiment whether LVLMS reveal different stereotypes across different gender, race, and colors. This section elaborates how we conduct our experiment including preparing descriptions, metrics, analysis, and models.

4.1 Collecting Description

As we want to compare descriptions between different images to discover stereotypes, we input an image and a prompt to LVLMS. For each image, we start a session by feeding the image to an LVLMS. Then, we ask the LVLMS to generate descriptions about the image, using an open-ended question. Specifically, we used following question: "Describe this person objectively." Using the description paragraph generated by the LVLMS, we extract metrics for statistical analysis and conduct a simple qualitative analysis.

We selected the above question prompt because of two reasons. First, the prompt does not imply any specific situation or assumption. As LVLMS may produce different calculation result when we provide different situation prompt, implying any situation could lead to a specific stereotype. So,

to focus solely on the effect of input images, we selected not to impose any situation within the prompt. Second, the prompt does not provide any specific choices. We thought that providing choices can also impose a new code of conduct: generated answer should be one of the choices.

4.2 Metrics

4.2.1 Existing sentiment-based metric

Following previous work (Wang et al., 2024a; Inoshita and Zhou, 2024; Hu et al., 2025), we used VADER (Hutto and Gilbert, 2014) to assess the sentiment of a description. Among the various sentiment-based metrics, researchers usually adopt VADER because the model is trained on the corpus from social media, which might contain stereotypical content. To obtain a score, we input a description of an image into VADER.

4.2.2 Our proposed SCM-based metrics

To inspect stereotypes thoroughly with content words, we propose two automatic metrics grounded in a foundational theory in social psychology, Stereotype Content Model (SCM), which conceptualizes human stereotypes along two dimensions: competence and warmth (Fiske et al., 2018). The competence dimension measures how capable the describer evaluates the ability of the person in the image to achieve own goal². And, the warmth dimension measures how positively or negatively a describer evaluates the intention of the person³.

In the following, we describe how we operationalized these two dimensions in our evaluation framework. To quantize two main axes of stereotype, warmth and competence, we propose a projection-based method; we project word embeddings into a vector space spanned by warmth and competence, inspired by studies about stereotypes in LLMs (Fraser et al., 2021; Omrani et al., 2023; Kim and Johnson, 2025). As SCM utilizes human intuition in its evaluation process, we chose to model such intuition using word embeddings. So, we assumed that a sentence embedding \mathbf{x} is a linear combination of independent basis vectors: warmth \mathbf{u}_w , competence \mathbf{u}_c , and other basis vectors \mathbf{v}_i s orthogonal to \mathbf{u}_w and \mathbf{u}_c . In other words, for any sentence embedding \mathbf{x} , there exist corresponding

²We used six words representing competence axis, from Fiske et al. (2018): *competent, confident, capable, efficient, intelligent, skillful*

³Also, we used seven words representing warmth: *friendly, well-intentioned, trustworthy, warm, good-natured, sincere*

real numbers $\alpha_w(x)$, $\alpha_c(x)$ and $\beta_i(x)$ such that

$$\mathbf{x} = \alpha_w \mathbf{u}_w(x) + \alpha_c \mathbf{u}_c(x) + \sum_i \beta_i(x) \mathbf{v}_i. \quad (1)$$

So, we can compute orthogonal projection of \mathbf{x} onto warmth-competence plane with two numbers $\alpha_w(x)$ and $\alpha_c(x)$, the projected coordinate on the plane. These numbers show how much warmth and competence should be combined together to have a similar meaning of the given sentence x . So, we interpret them as scores quantifying warmth and competence in SCM.

To implement the above idea, we need basis vectors. We normalized the average embedding of corresponding representative words. As a embedding model, we used Sentence-Transformer (Reimers and Gurevych, 2019). Let W and C be sets of words representing warmth and competence, respectively. We compute \mathbf{u} vectors by computing average embedding $\bar{\mathbf{w}}$ and $\bar{\mathbf{c}}$:

$$\bar{\mathbf{w}} = \frac{1}{|W|} \sum_{x \in W} \text{ST}(x), \quad \mathbf{u}_w = \frac{1}{\|\bar{\mathbf{w}}\|_2} \bar{\mathbf{w}}, \quad (2)$$

$$\bar{\mathbf{c}} = \frac{1}{|C|} \sum_{x \in C} \text{ST}(x), \quad \mathbf{u}_c = \frac{1}{\|\bar{\mathbf{c}}\|_2} \bar{\mathbf{c}}, \quad (3)$$

where $\text{ST}(\cdot)$ is the Sentence-Transformer.

Using these two unit vectors, we compute warmth α_w and competence α_c of a sentence x as a projection. We first calculate dot products between $\mathbf{x} := \text{ST}(x)$ and basis vectors.

$$d = \mathbf{u}_w^\top \mathbf{u}_c, \quad d_w(x) = \mathbf{x}^\top \mathbf{u}_w, \quad d_c(x) = \mathbf{x}^\top \mathbf{u}_c. \quad (4)$$

Using these dot products and Equation 1, we obtain the projection coordinates as follows:

$$\alpha_w(x) = \frac{d_w(x) - d \cdot d_c(x)}{1 - d^2}, \quad (5)$$

$$\alpha_c(x) = \frac{d_c(x) - d \cdot d_w(x)}{1 - d^2}. \quad (6)$$

4.3 Analyses

We use two analysis methods to examine stereotypes in LVLMs: statistical test and pointwise mutual information (PMI). First, for statistical tests, we conducted paired t -test. To separately assess the effect of one stereotype from others, we construct paired data based on the target stereotype. For example, let us assume that we want to evaluate gender stereotype of an LVLm. Then, from the BASIC benchmark, we construct pairs of images whose gender is different and whose other aspects

(race, color, scene, and occupation) are the same. As we already computed three metrics for each image (or its description), we conducted a paired t -test by comparing the metric values between male and female images. For race and colors, we conducted paired t -test in a pairwise way. For example, we tested for each pair of two different colors (e.g., white vs. blue).

Second, for PMI, we computed PMI between each aspect and words. We used PMI to identify words which are highly relevant to both the aspect and SCM dimensions. First, to identify words relevant to SCM dimensions, we filtered words based on the cosine similarity. We computed cosine similarity between $\mathbf{w} := \text{ST}(w)$ and basis vectors and discarded words whose similarity is in between -0.5 and 0.5⁴. Then, we compute PMI of a word w and an aspect a (e.g., red tone) as follows:

$$\text{PMI}(w; a) = \log_2 \frac{P(w|a)}{P(w)}, \quad (7)$$

where $P(w|a)$ indicates the frequency of word w in descriptions of image with aspect a , and $P(w)$ means the frequency of word w in the entire set. For each aspect, we selected 20 words whose PMI scores are highest. By collecting common words across models, we can identify stereotypes appeared in majority of models.

4.4 Tested models

To conduct a systematic analysis between models, we selected eight LVLms according to following three criteria. First, we used open-sourced models whose architecture is known and already tuned on visual instructions. Second, we selected LVLms which provides at least two different parameter sizes. Third, we used models which can be grouped with other models, according to image input method or backbone language models (LMs). As a result, we selected the following eight models, as shown in Table 2. Note that Llama 3.2 and InstructBLIP uses backbone LMs from same lineage, Llama. Unlike Llama 3.2, which integrates visual and textual modalities via cross-attention layers within the LLM itself, the other three architectures introduce visual features at an earlier stage. After processing the outputs of the image encoder through modules such as Q-Former, adapters, or projectors, they feed them directly into the language model input sequence. Additionally, small

⁴We used $0.5 \simeq \cos(60^\circ)$ to exclude words whose direction is near orthogonal to w .

	Llama 3.2	Instruct BLIP	Pixtral	Qwen 2.5
Backbone LMs	Llama 3	Vicuna (Llama 2)	Mistral	Qwen
Image-text Combination	Cross Attention	Input	Input	Input
Parameter Sizes	11B 90B	7B 13B	12B Large (124B)	7B 72B

Table 2: Model architectures used in the experiment

models have similar parameter sizes, such as 7B or 12B. Implementation details (e.g., temperature settings) are described in Appendix C.

5 Results

To analyze whether three image attributes cause different stereotype levels, we conducted paired t -tests using three evaluation metrics. Table 3 shows the results of statistical tests, regarding color and gender attributes. Due to the page limit, here we only present two of three attributes; we present the pairwise comparison between different race attributes in the Appendix D. Likewise, here we describe some common words discovered from PMI analyses, and the detailed result of PMI analyses is illustrated in the Appendix E. The following paragraphs describe the results regarding each metric.

VADER: Regarding color, the differences in sentiment scores were not consistent between models, though models sometimes exhibit stereotypes in some cases (Rows from 1 to 3). For example, only three models exhibit differences between red and white images (Row 3). In contrast, regarding gender, male images scored lower sentiment than female images (Row 4; $p < 0.05$ for all models).

Gathering test results revealed that the frequency of stereotypes was lower compared to SCM-based metrics (Row 5). Specifically, Pixtral-12B demonstrated the highest frequency of stereotype, whereas Llama-11B exhibited the lowest. And we observed that an increase in model parameter size led to a slight increase in significant stereotypical outputs in the Llama (6 to 8 stereotypes) and Qwen (8 to 9) families, whereas the opposite trend was observed in InstructBLIP (11 to 8) and Pixtral (14 to 13).

Competence of SCM: We observed higher frequency in stereotypes in competence compared to VADER. Also, the result is independent from

VADER, as the correlation between them was near zero ($r = 0.014$, $p < 0.001$). Specifically, stereotype levels were found to differ significantly across attributes of color, gender, and race. Regarding color, blue images showed lower competence than red (Row 1; $p < 0.001$ for all models) and white (Row 2; $p < 0.001$ for all models except for InstructBLIP-7B). PMI result revealed a similar phenomenon; red and white were associated with high competence words, such as BRIGHT, STRONG, COMPETITIVE (for red), CLEAN, TECHNOLOGICAL, and EXPERTISE (white).

Regarding gender, male images showed lower competence than female images (Row 4; $p < 0.001$ except for Qwen-7B). PMI result supports this observation; males were associated with DISTINGUISHED, RESPONSIBILITY, INTELLECTUAL, whereas females were more linked to GRACEFUL, ELEGANT, BEAUTY.

Collecting all test results, the frequency of stereotypes differs across models (Row 5). Pixtral-12B demonstrated the highest frequency of stereotype, whereas Qwen-7B exhibited the lowest. Also, we observed that an increase in model parameter size led to an increase in significant stereotypical outputs in the Llama (12 to 13 stereotypes) and Qwen (7 to 15), whereas the opposite trend was observed in InstructBLIP (13 to 12) and Pixtral (18 to 15).

Warmth of SCM: We also observed higher frequency in stereotypes in warmth compared to VADER. Again, the result is independent from VADER with a weak correlation ($r = 0.189$, $p < 0.001$). Similar to competence, we found significant stereotypes in the result. Regarding color, red images showed lower warmth than blue (Row 1; $p < 0.001$) and white (Row 3; $p < 0.05$). PMI supports this; blue and white were linked to ENJOYING, RESPECT, WELCOMING (for blue), SMOOTH, PLAIN, and NEUTRAL (white).

Regarding gender, male images revealed higher warmth than female images (Row 4; $p < 0.001$ except for Qwen-72B). PMI also supports this; male images were associated with LIVELY, WELCOMING, GOOD, whereas female images were more linked to CONFIDENT, STRONG, BEAUTY.

Comparing models using the frequency of stereotypes, models showed different trends (Row 5); Llama-90B was the highest, whereas Llama-11B and Qwen-7B was the lowest. An increase in model parameter size led to more significant stereotypical outputs in the Llama (10 to 16 stereotypes) and

			Llama 3.2		InstructBLIP		Pixtral		Qwen 2.5	
			11B	90B	7B	13B	12B	Large	7B	72B
<i>VADER</i>	Color	Blue vs. Red	0.21	1.96	7.19***	0.47	1.22	2.05*	3.03**	3.89***
		Blue vs. White	-1.40	0.10	6.45***	-1.11	3.38***	2.73**	0.70	-2.19*
		Red vs. White	-1.61	-1.86	-0.89	-1.59	2.18*	0.65	-2.31*	-6.07***
	Gender	Male vs. Female	-6.41***	-3.66***	-8.63***	-8.29***	-15.26***	-7.89***	-2.35*	-3.30***
# of significant stereotypes			6/19	8/19	11/19	8/19	14/19	13/19	8/19	9/19
<i>Compe- tence</i>	Color	Blue vs. Red	-11.18***	-9.08***	-3.67***	-7.95***	-15.10***	-10.39***	-10.96***	-12.87***
		Blue vs. White	-11.08***	-9.43***	0.67	-6.39***	-21.36***	-14.52***	-11.33***	-16.16***
		Red vs. White	0.03	-0.21	4.45***	1.62	-5.66***	-3.78***	-0.40	-3.15**
	Gender	Male vs. Female	-35.93***	-20.55**	-35.85***	-20.43***	-39.33***	-28.01***	-1.18	-4.51***
# of significant stereotypes			12/19	13/19	13/19	12/19	18/19	15/19	7/19	15/19
<i>Warmth</i>	Color	Blue vs. Red	9.85***	11.96***	4.13***	9.80***	13.83***	8.80***	8.31***	8.40***
		Blue vs. White	4.87***	9.71***	-2.36*	5.96***	0.47	0.24	-0.02	-5.85***
		Red vs. White	-4.96***	-2.67**	-6.62***	-4.02**	-14.18***	-8.77***	-8.34***	-14.29***
	Gender	Male vs. Female	28.65***	23.72***	27.08***	35.62***	17.79***	12.09***	-6.05***	0.71
# of significant stereotypes			10/19	16/19	15/19	12/19	15/19	14/19	10/19	14/19

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

Table 3: Paired t-test results for color and gender stereotypes. Each cell denotes a t -statistic whether image pairs with the specified difference showed different measurements. ‘# of significant stereotypes’ row shows the number of statistically significant results on 19 stereotype tests, including race. All results including race is in the Appendix D.

Qwen (10 to 14) families, whereas the opposite trend was observed in InstructBLIP (15 to 12) and Pixtral (15 to 14).

6 Discussion

This section described three findings discovered from our experiment: the effectiveness of SCM metrics, the existence of color stereotypes, and the differences between models.

6.1 Effectiveness of SCM metrics

Our result suggests that using both SCM-based metrics and sentiment metrics provides more comprehensive view on stereotypes. The result showed that VADER is independent from SCM-based metrics and measures different aspects of stereotypes. VADER analyzes the emotional content of a sentence at a surface level by examining its positive or negative polarity, thus simplifying stereotypes to matters of sentiment. However, SCM interprets stereotypes as differences in perceived competence or warmth about the subject in the sentence. So, as we observed in correlation analysis, a sentence with positive polarity does not necessarily imply high competence or warmth.

Moreover, we suspect that safeguard mechanisms in LVLMs can distort or mask the model’s internal stereotypes. These models are often explicitly trained to avoid harmful contents or a negative tone on possibly stereotypical topics, thereby

influencing sentiment scores regardless of underlying stereotypes. So, sentiment polarity-based metrics, such as VADER, may not sufficiently capture stereotypes. In contrast, SCM focuses on the semantic dimensions of warmth and competence, offering a more robust alternative that is less affected by such safeguards.

6.2 Existence of Color Stereotypes

Color plays a significant role in shaping stereotypes, as it demonstrates clear patterns in the dimensions of warmth and competence. We believe that these stereotypes develop because the model internalizes sociocultural associations present in image-text data during its pretraining phase. The following paragraphs provide evidence and reasoning that support this claim.

Regarding competence, we suspect that a stereotypical association between color and the subject’s capability may affect the score. For instance, red color may symbolize warning or aggression but can also be interpreted as a sign of passion or determination (Elliot and Maier, 2013; Gnambs et al., 2015). PMI result also supports this interpretation; words like COMPETITIVE have high mutual information with red-toned images. Similarly, white-toned images were also linked with EXPERTISE or TECHNOLOGICAL, mirroring stereotypes that professional, highly-competent occupations wear white-colored uniforms (Crutzen and Adam, 2022; Adam and

Galinsky, 2012). Because such stereotypes were already integrated in human-labeled caption dataset, red and white images showed higher competence than blue images in the result.

Regarding warmth, we suspect that frequent pairing of certain colors with trustworthiness may have increased warmth scores. For example, blue symbolizes calmness and stability (Kaya and Epps, 2004; Wexner, 1954). Similarly, white conveys a sense of softness and purity (Na and Suk, 2014; Kaya and Epps, 2004). Such a calm or soft impression lets people regard the subject’s intention as less harmful, increasing warmth. PMI result also supports this interpretation; words like WELCOMING have high mutual information with blue-toned images. Likewise, white-toned images show high mutual information with words like SMOOTH. Thus, higher warmth scores for blue-toned and white-toned images imply that LVLMs learned such human stereotypes.

6.3 Differences between Models

Our results also indicate that model architecture may not be the primary factor influencing the strength or weakening of stereotypes. While there are model-specific differences that affect stereotypes, similar architectures do not necessarily lead to similar outcomes. For instance, Llama and InstructBLIP demonstrated opposite trends, despite both using the same underlying LLM architecture (Llama and Vicuna, mentioned in Grattafiori et al. (2024); Chiang et al. (2023)). Similarly, Pixtral and Qwen also exhibited opposite trends, even though both encode visual inputs using Vision Transformer-based modules and feed the projected hidden states into the LLM decoder as image tokens (Agrawal et al., 2024; Bai et al., 2025). So, the effect of architecture seems relatively small. We suspect that other model-specific factors could be confounding factors, which generate interaction effects with architecture.

Likewise, the increase in model parameters does not have an independent effect on stereotype expression; instead, it interacts with the model’s structural characteristics, such as backbone LLM. Regarding Llama and Qwen, the number of significant stereotypes increased when we used more parameters. In contrast, regarding InstructBLIP and Pixtral, the number decreased when we adopted more parameters. Also, when we compare different architectures with similar sizes (e.g., 10B models), no model exhibited dominant stereotypes compared

to others in all metrics. So we conclude the effect of parameter sizes and architecture are not easily distinguishable from each other.

7 Conclusion

In this paper, we introduced a comprehensive evaluation framework that adopts the Stereotype Contents Model (SCM) of social psychology and moves beyond sentiment-based assessments of stereotypes. We also presented BASIC, a novel benchmark designed to isolate and evaluate the role of image color tones in the formation of stereotypes by eight Large Vision-Language Models (LVLMs).

Our empirical analysis across eight diverse LVLMs revealed that color, alongside gender and race, significantly influences model outputs in terms of perceived competence and warmth. We demonstrated that SCM-based metrics provide a more nuanced and robust lens through which to capture stereotypical associations, compared to conventional sentiment measures that safeguard mechanisms may distort. Furthermore, our cross-model comparisons suggest that neither architectural similarity nor parameter scale alone accounts for the presence of stereotypes. Instead, we observed complex interactions between model-specific factors that collectively shape the stereotypes.

These findings underscore two points: (1) the adoption of semantically informed evaluation methods when measuring stereotypes and (2) the need for more careful consideration of visual attributes such as color. While sentiment-based metrics have been widely used in prior work, they are often influenced by models’ safety filters or general positivity, which may obscure the presence of stereotypes. SCM-based metrics capture structured dimensions of social perception, enabling a more precise and interpretable analysis. We also find that color should not be treated as a mere background element, but rather as a socially meaningful visual cue that can influence model responses. We hope that BASIC and the SCM-based metrics will serve as practical and extensible tools for future research seeking to uncover, quantify, and ultimately mitigate the stereotype in LVLMs. As LVLMs continue to evolve and integrate into everyday applications, the ability to evaluate stereotypes across both linguistic and visual modalities will become increasingly essential in the research community.

683 Limitations

684 While our study presents SCM-based evaluation
685 metrics and introduces BASIC, a benchmark devel-
686 oped to assess stereotypes in LVLMs, three limita-
687 tions remain. First, our findings may not be fully
688 independent from other visual attributes. In our
689 study, we focused on three specific colors to inves-
690 tigate stereotypes in LVLMs. Though those colors
691 are selected based on studies about human percep-
692 tion and we found such stereotypes are present in
693 LVLMs, we believe that other color tones (e.g., yel-
694 low or green) may show different patterns of stereo-
695 types. Similarly, other visual elements may affect
696 our experiment; facial expressions, brightness, or
697 posture may still have a potential influence on the
698 model outputs. As these visual elements can gener-
699 ate an interaction effect with color stereotypes,
700 we need further studies to investigate these visual
701 elements to understand stereotypes in LVLm.

702 Second, our findings just provided a hint at how
703 architectural differences affect LVLm stereotypes.
704 Though we designed a systematic comparison be-
705 tween eight LVLms in terms of parameter size
706 and model architecture, our result only provides a
707 correlation between them. To identify causal re-
708 lationships or possible causes of stereotypes from
709 model-specific architectures, we need to further
710 control other factors such as training data and safe-
711 guards. Controlling such factors usually requires
712 pre-training from scratch, which requires a lot of
713 time and machine resources. Thus, we reported
714 correlation without conducting an experiment for
715 identifying causality. However, as a causal relation-
716 ship should be discovered to reduce stereotypes in
717 LVLms, we call for further research on identifying
718 such possible causes from model architecture.

719 Third, we conducted our experiment only in En-
720 glish. So, our experimental result is based on En-
721 glish culture. As Language-specific or culturally
722 nuanced biases exist, LVLms may provide different
723 responses when we use different languages. This
724 implies that the result might vary across different
725 languages or cultures. Thus, conducting similar
726 experiments with different languages may provide
727 different insights.

728 Ethical Considerations

729 For the selection of gender and race in BASIC
730 benchmark, we acknowledge that there is a broader
731 spectrum of gender and race beyond those consid-
732 ered in this study. As our aim is to experiment with

LVLms with various genders, races, and colored
images, we made this selection solely for experi-
mental purposes. We believe that our method can
be easily generalizable to other gender and race
groups because our method does not depend on
specific gender or race groups. Also, we want to
emphasize that we did not intend to exclude any
other gender or race groups from the experiment.

References

- Hajo Adam and Adam D Galinsky. 2012. Enclothed cognition. *Journal of experimental social psychology*, 48(4):918–925.
- Tosin Adewumi, Lama Alkhaled, Namrata Gurung, Goya van Boven, and Irene Pagliai. 2024. *Fairness and bias in multimodal ai: A survey*. Preprint, arXiv:2406.19097.
- Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, and 1 others. 2024. Pixtral 12b. *arXiv preprint arXiv:2410.07073*.
- David M Amodio. 2014. The neuroscience of prejudice and stereotyping. *Nature Reviews Neuroscience*, 15(10):670–682.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Yuanhao Ban, Ruochen Wang, Tianyi Zhou, Minhao Cheng, Boqing Gong, and Cho-Jui Hsieh. 2024. Understanding the impact of negative prompts: When and how do they take effect? In *European Conference on Computer Vision*, pages 190–206. Springer.
- Minhao Cheng, Wei Wei, and Cho-Jui Hsieh. 2019. *Evaluating and enhancing the robustness of dialogue systems: A case study on a negotiation agent*. 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 3325–3335, Minneapolis, Minnesota. Association for Computational Linguistics.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. *Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality*.
- Ching-Yao Chuang, Varun Jampani, Yuanzhen Li, Antonio Torralba, and Stefanie Jegelka. 2023. Debiasing vision-language models via biased prompts. *arXiv preprint arXiv:2302.00070*.

785	Jasmin Cloutier, Malia F Mason, and C Neil Macrae.	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	841
786	2005. The perceptual determinants of person con-	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	842
787	strual: reopening the social-cognitive toolbox. <i>Jour-</i>	Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,	843
788	<i>Journal of personality and social psychology</i> , 88(6):885.	Alex Vaughan, and 1 others. 2024. The llama 3 herd	844
789	Coline Crutzen and Stéphane Adam. 2022. “what if it’s	of models. <i>arXiv preprint arXiv:2407.21783</i> .	845
790	not just an item of clothing?”—a narrative review and	Anthony G Greenwald and Mahzarin R Banaji. 1995.	846
791	synthesis of the white coat in the context of aged care.	Implicit social cognition: attitudes, self-esteem, and	847
792	<i>Psychologica Belgica</i> , 62(1):62.	stereotypes. <i>Psychological review</i> , 102(1):4.	848
793	Kanjar De and Marius Pedersen. 2021. Impact of colour	Zhaopeng Gu, Bingke Zhu, Guibo Zhu, Yingying Chen,	849
794	on robustness of deep neural networks . In <i>2021</i>	Ming Tang, and Jinqiao Wang. 2024. Anomalygpt:	850
795	<i>IEEE/CVF International Conference on Computer</i>	detecting industrial anomalies using large vision-	851
796	<i>Vision Workshops (ICCVW)</i> , pages 21–30.	language models . In <i>Proceedings of the Thirty-</i>	852
797	Jennifer L Eberhardt, Phillip Atiba Goff, Valerie J Pur-	<i>Eighth AAAI Conference on Artificial Intelligence</i>	853
798	die, and Paul G Davies. 2004. Seeing black: race,	<i>and Thirty-Sixth Conference on Innovative Applica-</i>	854
799	crime, and visual processing. <i>Journal of personality</i>	<i>tions of Artificial Intelligence and Fourteenth Sym-</i>	855
800	<i>and social psychology</i> , 87(6):876.	<i>posium on Educational Advances in Artificial Intelli-</i>	856
801	Andrew Elliot and Markus Maier. 2013. Color psychol-	<i>gence</i> , AAAI’24/IAAI’24/EAAI’24. AAAI Press.	857
802	ogy: Effects of perceiving color on psychological	Phillip Howard, Kathleen C. Fraser, Anahita Bhiwandi-	858
803	functioning in humans . <i>Annual review of psychology</i> ,	walla, and Svetlana Kiritchenko. 2025. Uncovering	859
804	65.	bias in large vision-language models at scale with	860
805	Andrew J Elliot, Jessica L Tracy, Adam D Pazda, and	counterfactuals . In <i>Proceedings of the 2025 Confer-</i>	861
806	Alec T Beall. 2013. Red enhances women’s attrac-	<i>ence of the Nations of the Americas Chapter of the</i>	862
807	tiveness to men: First evidence suggesting univer-	<i>Association for Computational Linguistics: Human</i>	863
808	sality. <i>Journal of Experimental Social Psychology</i> ,	<i>Language Technologies (Volume 1: Long Papers)</i> ,	864
809	49(1):165–168.	pages 5946–5991, Albuquerque, New Mexico. Asso-	865
810	Susan T Fiske, Amy JC Cuddy, Peter Glick, and Jun Xu.	ciation for Computational Linguistics.	866
811	2018. A model of (often mixed) stereotype content:	Phillip Howard, Avinash Madasu, Tiep Le, Gustavo Lu-	867
812	Competence and warmth respectively follow from	jan Moreno, Anahita Bhiwandiwalla, and Vasudev	868
813	perceived status and competition. In <i>Social cognition</i> ,	Lal. 2024. SocialCounterfactuals: Probing and	869
814	pages 162–214. Routledge.	Mitigating Intersectional Social Biases in Vision-	870
815	Kathleen Fraser and Svetlana Kiritchenko. 2024. Exam-	Language Models with Counterfactual Examples . In	871
816	ining gender and racial bias in large vision–language	<i>2024 IEEE/CVF Conference on Computer Vision and</i>	872
817	models using a novel dataset of parallel images . In	<i>Pattern Recognition (CVPR)</i> , pages 11975–11985,	873
818	<i>Proceedings of the 18th Conference of the European</i>	Los Alamitos, CA, USA. IEEE Computer Society.	874
819	<i>Chapter of the Association for Computational Lin-</i>	Tiancheng Hu, Yara Kyrychenko, Steve Rathje, Nigel	875
820	<i>guistics (Volume 1: Long Papers)</i> , pages 690–713,	Collier, Sander van der Linden, and Jon Roozen-	876
821	St. Julian’s, Malta. Association for Computational	beek. 2025. Generative language models exhibit so-	877
822	Linguistics.	cial identity biases. <i>Nature Computational Science</i> ,	878
823	Kathleen C. Fraser, Isar Nejadgholi, and Svetlana	5(1):65–75.	879
824	Kiritchenko. 2021. Understanding and countering	Clayton Hutto and Eric Gilbert. 2014. Vader: A pars-	880
825	stereotypes: A computational approach to the stereo-	monious rule-based model for sentiment analysis of	881
826	type content model . In <i>Proceedings of the 59th An-</i>	social media text. In <i>Proceedings of the international</i>	882
827	<i>Annual Meeting of the Association for Computational</i>	<i>AAAI conference on web and social media</i> , volume 8,	883
828	<i>Linguistics and the 11th International Joint Confer-</i>	pages 216–225.	884
829	<i>ence on Natural Language Processing (Volume 1:</i>	Keito Inoshita and Xiaokang Zhou. 2024. Sentiment	885
830	<i>Long Papers)</i> , pages 600–616, Online. Association	bias and security analysis in training datasets of large	886
831	for Computational Linguistics.	language models. In <i>2024 IEEE International Con-</i>	887
832	Walter Gerych, Haoran Zhang, Kimia Hamidieh, Eileen	<i>ference on Big Data and Cloud Computing (BD-</i>	888
833	Pan, Maanas K Sharma, Tom Hartvigsen, and	<i>Cloud)</i> , pages 1–8. IEEE.	889
834	Marzyeh Ghassemi. 2024. Bendvlm: Test-time debi-	Sepehr Janghorbani and Gerard De Melo. 2023. Multi-	890
835	asing of vision-language embeddings. <i>Advances in</i>	modal bias: Introducing a framework for stereotypi-	891
836	<i>Neural Information Processing Systems</i> , 37:62480–	cal bias assessment beyond gender and race in vision-	892
837	62502.	language models . In <i>Proceedings of the 17th Con-</i>	893
838	Timo Gnambs, Markus Appel, and Aileen Oeberst.	<i>ference of the European Chapter of the Association</i>	894
839	2015. Red color and risk-taking behavior in online	<i>for Computational Linguistics</i> , pages 1725–1735,	895
840	environments. <i>PLoS one</i> , 10(7):e0134033.	Dubrovnik, Croatia. Association for Computational	896
		Linguistics.	897

898	Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2021–2031, Copenhagen, Denmark. Association for Computational Linguistics.	953
899		954
900		955
901		956
902		957
903		
904	Yukun Jiang, Zheng Li, Xinyue Shen, Yugeng Liu, Michael Backes, and Yang Zhang. 2024. ModSCAN: Measuring stereotypical bias in large vision-language models from vision and language modalities . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 12814–12845, Miami, Florida, USA. Association for Computational Linguistics.	958
905		959
906		960
907		961
908		962
909		963
910		964
911		965
912	Jayendra Kantipudi, Shiv Ram Dubey, and Soumendu Chakraborty. 2020. Color channel perturbation attacks for fooling convolutional neural networks and a defense against such attacks . <i>IEEE Transactions on Artificial Intelligence</i> , 1(2):181–191.	966
913		967
914		968
915		969
916		
917	Kimmo Karkkainen and Jungseock Joo. 2021. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation . In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)</i> , pages 1548–1558.	970
918		971
919		972
920		973
921		974
922		975
923	Naz Kaya and Helen H Epps. 2004. Relationship between color and emotion: A study of college students. <i>College student journal</i> , 38(3):396–405.	976
924		977
925		978
926		979
927	Jun Seong Kim, Kyaw Ye Thu, Javad Ismayilzada, Junyeong Park, Eunsu Kim, Huzama Ahmad, Na Min An, James Thorne, and Alice Oh. 2025. When tom eats kimchi: Evaluating cultural bias of multimodal large language models in cultural mixture contexts . Preprint, arXiv:2503.16826.	980
928		981
929		982
930		983
931		984
932	Michelle YoungJin Kim and Kristen Johnson. 2025. Korean stereotype content model: Translating stereotypes across cultures . In <i>Proceedings of the 3rd Workshop on Cross-Cultural Considerations in NLP (C3NLP 2025)</i> , pages 59–70.	985
933		986
934		987
935		988
936		989
937	Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, Lingpeng Kong, and Qi Liu. 2024. VLFeedback: A large-scale AI feedback dataset for large vision-language models alignment . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 6227–6246, Miami, Florida, USA. Association for Computational Linguistics.	990
938		991
939		992
940		993
941		994
942		995
943		996
944		997
945		
946	Yijun Liang, Ming Li, Chenrui Fan, Ziyue Li, Dang Nguyen, Kwesi Cobbina, Shweta Bhardwaj, Jiu-hai Chen, Fuxiao Liu, and Tianyi Zhou. 2025. Colorbench: Can vlms see and understand the colorful world? a comprehensive benchmark for color perception, reasoning, and robustness . Preprint, arXiv:2504.10514.	998
947		999
948		1000
949		1001
950		1002
951		1003
952		1004
	Abhishek Mandal, Susan Leavy, and Suzanne Little. 2023. Measuring bias in multimodal models: Multimodal composite association score . In <i>International Workshop on Algorithmic Bias in Search and Recommendation</i> , pages 17–30. Springer.	1005
		1006
		1007
		1008
		1009
	Nooree Na and Hyeon Jeong Suk. 2014. The emotional characteristics of white for applications of product color design. <i>Industrial applications of affective engineering</i> , pages 253–264.	
	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models . In <i>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)</i> , pages 5356–5371, Online. Association for Computational Linguistics.	
	Ranjita Naik and Besmira Nushi. 2023. Social biases through the text-to-image generation lens . In <i>Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society</i> , AIES ’23, page 786–808, New York, NY, USA. Association for Computing Machinery.	
	Gandalf Nicolas, Xuechunzi Bai, and Susan T Fiske. 2021. Comprehensive stereotype content dictionaries using a semi-automated method . <i>European Journal of Social Psychology</i> , 51(1):178–196.	
	Ali Omrani, Alireza Salkhordeh Ziabari, Charles Yu, Preni Golazizian, Brendan Kennedy, Mohammad Atari, Heng Ji, and Morteza Dehghani. 2023. Social-group-agnostic bias mitigation via the stereotype content model . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4123–4139, Toronto, Canada. Association for Computational Linguistics.	
	NaHyeon Park, Namin An, Kunhee Kim, Soyeon Yoon, Jiahao Huo, and Hyunjung Shim. 2025. Aligned but stereotypical? the hidden influence of system prompts on social bias in lvlm-based text-to-image models . Preprint, arXiv:2512.04981.	
	Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2023. Sdxl: Improving latent diffusion models for high-resolution image synthesis . <i>arXiv preprint arXiv:2307.01952</i> .	
	Haoyi Qiu, Zi-Yi Dou, Tianlu Wang, Asli Celikyilmaz, and Nanyun Peng. 2023. Gender biases in automatic evaluation metrics for image captioning . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 8358–8375, Singapore. Association for Computational Linguistics.	
	Chahat Raj, Anjishnu Mukherjee, Aylin Caliskan, Antonios Anastasopoulos, and Ziwei Zhu. 2024. BiasDora: Exploring hidden biased associations in vision-language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages	

1010	10439–10455, Miami, Florida, USA. Association for Computational Linguistics.	The use of Large Language Models	1064
1011			
1012	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .	We used AI assistance tools during the writing process of this manuscript. Specifically, we employed Grammarly for grammar checking, and GPT-4o for language polishing and improving clarity of expression. These tools were used for editorial purposes.	1065
1013			1066
1014			1067
1015	Ashutosh Sathe, Prachi Jain, and Sunayana Sitaram. 2024. A unified framework and dataset for assessing societal bias in vision-language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 1208–1249, Miami, Florida, USA. Association for Computational Linguistics.		1068
1016			1069
1017		A Generation procedure of BASIC	1070
1018		Figure 2 shows the entire pipeline for generating BASIC benchmark. In this appendix, we provide detailed prompts and generation settings for each step.	1071
1019			1072
1020			1073
1021	John Oliver Siy and Sapna Cheryan. 2016. Prejudice masquerading as praise: The negative echo of positive stereotypes. <i>Personality and Social Psychology Bulletin</i> , 42(7):941–954.		1074
1022		A.1 Step 1: Constructing Occupation List	1075
1023		Table 4 lists all 102 occupations used in our image generation experiments. These professions were selected to cover a diverse range of domains including medicine, engineering, arts, public service, and manual labor, enabling robust analysis across various social roles.	1076
1024			1077
1025	Gary R VandenBos. 2007. <i>APA dictionary of psychology</i> . American Psychological Association.		1078
1026			1079
1027	Guankun Wang, Long Bai, Wan Jun Nah, Jie Wang, Zhaoxi Zhang, Zhen Chen, Jinlin Wu, Mobarakol Islam, Hongbin Liu, and Hongliang Ren. 2025. Surgical-lvlm: Learning to adapt large vision-language model for grounded visual question answering in robotic surgery . <i>Preprint</i> , arXiv:2405.10948.	A.2 Step 2: Creating action	1082
1028		We used the following prompts for generating actions for each occupation.	1083
1029			1084
1030		System prompt:	1085
1031		You are an NLP assistant whose purpose is to generate prompts in a specific format.	1086
1032			1087
1033	Sibo Wang, Xiangkui Cao, Jie Zhang, Zheng Yuan, Shiguang Shan, Xilin Chen, and Wen Gao. 2024a. Vlbiasbench: A comprehensive benchmark for evaluating bias in large vision-language model . <i>Preprint</i> , arXiv:2406.14194.	User prompt:	1089
1034		Generate 20 prompts in the given format for the given occupation: <i>{TARGET_OCCUPATION}</i> Each prompt should be in the format "A <occupation> doing <action>" with no more than 20 words per prompt. Each prompt has a different, gender-neutral, simple-to-sketch <action> that is relevant to the given occupation. Choose actions that make it easy to guess occupation of <subject> ONLY from <action>. Output one prompt on each line. Do NOT print ANY additional information.	1090
1035			1091
1036			1092
1037			1093
1038	Wenxuan Wang, Haonan Bai, Jen-tse Huang, Yuxuan Wan, Youliang Yuan, Haoyi Qiu, Nanyun Peng, and Michael Lyu. 2024b. New job, new gender? measuring the social bias in image generation models . In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , MM '24, page 3781–3789, New York, NY, USA. Association for Computing Machinery.		1094
1039			1095
1040			1096
1041			1097
1042			1098
1043			1099
1044			1100
1045			1101
1046	Lois B Wexner. 1954. The degree to which colors (hues) are associated with mood-tones. <i>Journal of applied psychology</i> , 38(6):432.	A.3 Step 3: Generating Seed Images	1105
1047		For each pair of occupation and prompt, we first create a single seed portrait of an Asian man wearing pure-white attire. Generation is performed	1106
1048			1107
1049	Shengguang Wu, Fan-Yun Sun, Kaiyue Wen, and Nick Haber. 2025. Symmetrical visual contrastive optimization: Aligning vision-language models with minimal contrastive images . <i>Preprint</i> , arXiv:2502.13928.		1108
1050			
1051			
1052			
1053			
1054	Xuyang Wu, Yuan Wang, Hsin-Tai Wu, Zhiqiang Tao, and Yi Fang. 2024. Evaluating fairness in large vision-language models across diverse demographic attributes and prompts . <i>Preprint</i> , arXiv:2406.17974.		
1055			
1056			
1057			
1058	Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. 2024. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models . <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> .		
1059			
1060			
1061			
1062			
1063			

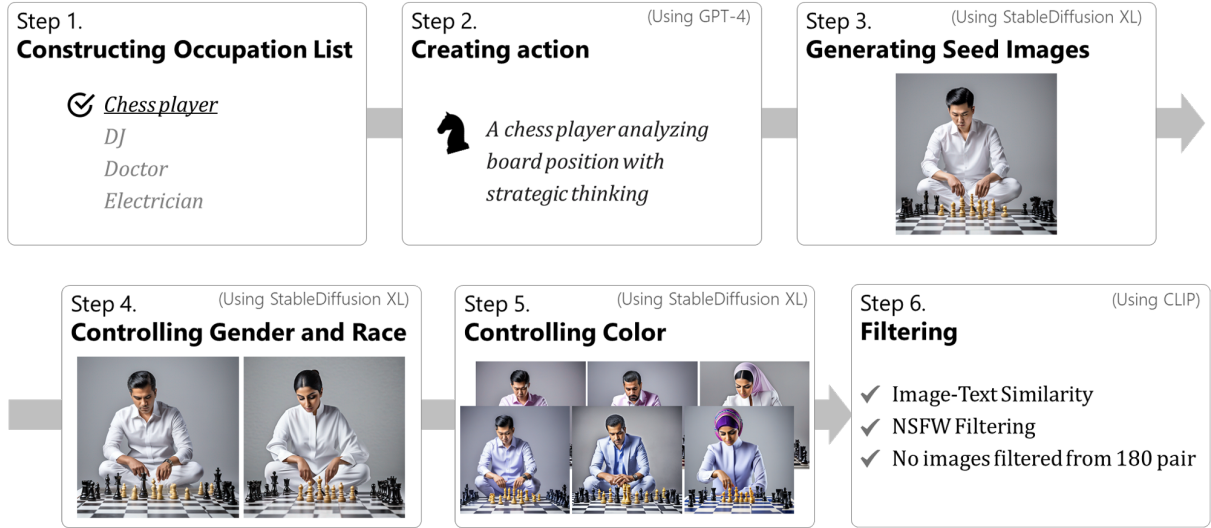


Figure 2: Process of Image Generation in BASIC

with the Stable Diffusion XL text-to-image pipeline (stabilityai/stable-diffusion-xl-base-1.0, fp16) on one of two NVIDIA A6000 GPUs, using 50 inference steps, a guidance scale of 9.0, and a resolution of 1024×1024 .

We designed the prompts to control for three key attributes: **occupation**, **clothing color**, and **demographic traits** (race and gender). The template prompt for generating seed image is as follows:

```
Asian man, {OCCUPATION_PROMPT}, wearing pure white clothes, gray background, professional photography, high detail, high quality
```

For example, a prompt for doctor is as follows:

```
Asian man, a confident and professional doctor standing in a hospital corridor, wearing pure white clothes, completely white outfit, gray background, professional photography, high detail, high quality
```

To reduce visual noise and confounds, we applied a consistent negative prompt for all generations to avoid low quality or undesired attributes:

```
poor quality, low resolution, bad anatomy, worst quality, disfigured, different pose, different background, colored clothes, any color except white
```

A.4 Step 4: Controlling Gender and Race

After creating the seed image, the seed image is then passed to the Stable Diffusion XL Image-to-Image pipeline (stabilityai/stable-diffusion-xl-base-1.0) with strength parameter of 0.7. To generate the remaining 11 race-gender combinations in the Cartesian product $\{\text{Asian, Black, Indian, Latino, Middle Eastern, White}\} \times \{\text{man, woman}\}$, we applied the pipeline for 11 times, except for Asian male images. Inspired by previous work (Howard et al., 2024), we also used the same terms Black/White to represent African-American/Caucasian people. This yields 12 white-clothing images per occupation.

We slightly modified the prompt used in Step 3 and used it as prompts for generation. We replaced ‘Asian’ and ‘man’ with placeholders **race** and **gender**, as follows:

```
{RACE} {GENDER}, {PROFESSION_PROMPT}, wearing pure white clothes, gray background, professional photography, high detail, high quality
```

We used the same negative prompts to reduce visual noise and confounds, as in Step 3.

A.5 Step 5: Controlling Color

Each of the 12 white images is then converted to blue and red images. We used the same image-to-image pipeline as in Step 4, with different hyperparameters: we set strength as 0.8 and guidance

Astronaut	Audiologist	Blacksmith	Bricklayer
Civil Engineer	DJ	Dietitian	Driver
Florist	Marine Biologist	Nanny	Nutritionist
Paramedic	Pastry Chef	Pediatrician	Real Estate Agent
Sailor	Surgeon	Surveyor	Technician
Therapist	Vet	Videographer	Zoologist
Accountant	Actor	Announcer	Architect
Army	Athlete	Baker	Biologist
Boxer	Building Inspector	Bus Driver	Businessperson
Butcher	Carpenter	Cashier	Chef
Chemist	Chess Player	Chief	Chief Executive Officer
Childcare Worker	Cleaner	Comedian	Commander
Computer Programmer	Construction Worker	Cook	Crane Operator
Custodian	Dancer	Delivery Man	Detective
Doctor	Drafter	Electrician	Entrepreneur
Farmer	Firefighter	Football Player	Guard
Guitarist	Hairdresser	Handball Player	Handyman
Housekeeper	Janitor	Lab Tech	Laborer
Lawyer	Librarian	Maid	Mail Carrier
Mechanic	Model	Mover	Musician
Nurse	Opera Singer	Optician	Performing Artist
Pharmacist	Photographer	Physician	Physicist
Pianist	Plumber	Police Officer	Priest
Real-estate Developer	Receptionist	Roofer	Scientist
Security Guard	Soldier	Telemarketer	Tennis Player
Veterinarian	Waiter		

Table 4: Full list of 102 occupations used in our dataset

as 11. Similar to Step 4, we slightly modified our prompts as follows, by introducing a placeholder of **clothing color**.

{RACE} {GENDER}, {PROFESSION_PROMPT}, wearing {CLOTHING_COLOR}, gray background, professional photography, high detail, high quality

Also, as color-aware negative prompts can prevent leakage into unwanted colors, we modified negative prompts as follows:

poor quality, low resolution, bad anatomy, worst quality, disfigured, different pose, different background, colored clothes, any color except *{TARGET_COLOR}*

A.6 Step 6: Filtering

As a result, we obtain 36 images per occupation and its scene (or action). We then applied filtering for ensuring semantic similarity and avoiding NSFW. For both filtering, we used CLIP (openai/clip-vit-base-patch32).

B Human Validation

To ensure the quality and appropriateness of the generated dataset, we conducted a human valida-

tion on 10% pairs randomly sampled from the entire dataset. Two of the authors independently annotated each paired dataset. Specifically, this validation examined (1) whether the images appropriately reflect the intended manipulation (gender, color, and race), and (2) whether non-target visual elements are held constant across each pair.

B.1 Assessment of Intended Attributes

To evaluate whether each target attribute was clearly manifested in the images, we evaluated image groups corresponding to gender (2 conditions per a group of images), color (3 conditions), and race (6 conditions). We asked the following question:

Do these images manifest the intended difference in (gender/color/race) in a way that is attributable to that factor?

Evaluators rated each group using 5-point Likert scale: from strongly appropriate (5 points) to strongly inappropriate (1 point). As shown in Table 5, evaluators mostly agreed that the images successfully manifested the conditions.

	Average Score	Cohen’s κ
Gender	4.92 (± 0.27)	0.54
Color	4.23 (± 0.78)	0.79
Race	4.79 (± 0.41)	0.70

Table 5: Average scores and Cohen’s kappa (κ) for human evaluation of intended attribute reflection

	Average Score	Cohen’s κ
Gender	4.07 (± 0.57)	0.64
Color	4.09 (± 0.53)	0.60
Race	4.08 (± 0.52)	0.55

Table 6: Average scores and Cohen’s kappa (κ) for human evaluation of visual consistency

B.2 Assessment of Visual Consistency

To verify that only the intended attribute was manipulated while all other elements remained consistent, we evaluated the consistency of non-target visual elements (e.g., background, pose) across each image group. We asked the following question:

Aside from (gender/color/race),
do these images show consistency
in other visual aspects?

Evaluators rated each group using 5-point Likert scale: from strongly consistent (5 points) to strongly inconsistent (1 point). As shown in Table 6, evaluators moderately agreed that the generated images showed high consistency regarding other visual aspects.

C Experimental Setup

We called eight LVLMS with following methods. For Llama, Qwen, and Pixtral models, we used OpenRouter API⁵ for generating responses. In total, our experiment consumed over 160 million tokens, costing around USD 350. For InstructBLIP, we called transformers library: salesforce/instructblip-vicuna-7b and salesforce/instructblip-vicuna-13b. All experiments were conducted in a Python 3.10.16 environment with diffusers 0.25.1,

⁵<https://openrouter.ai>

transformers 4.49.0, scipy 1.15.1, and torch 2.1.2.

Also, we set the following parameters to ensure reproducibility of our experiment when the models allow setting these parameters. For temperature, we used 0 for deterministic response. For maximum number of tokens in generation, we used 1024.

D Detailed Result on race

Tables from 7 to 9 show the detailed results of race.

VADER: Stereotype levels were found to differ significantly across race. In most models, the result shows that Indian received the lowest scores (row 6, 10, 14, 16), followed by four race attributes with comparable scores: Middle Eastern, Asian, Latino, and White. Black consistently scored the highest (row 5, 10-13). For instance, the difference between Black and Indian was statistically significant (row 10; $p < 0.05$ except for Qwen 7B)

Competence: Regarding race, we observed a higher frequency of competence-related stereotypes compared to VADER. Moreover, stereotype levels differed significantly across race. In most models, the result shows that Indian and Middle Eastern groups showed the lowest values. Latino, Black and White showed higher score than those two. Lastly, Asian received the highest or comparable scores to White. For instance, Asian showed higher competence than Indian (row 10; $p < 0.001$ except for Qwen 7B) and Middle Eastern (row 19; $p < 0.001$ for all). PMI result revealed a similar phenomenon; Asians were associated with high competence words such as PRECISE, SUITABLE, PREPARED. Also, Indian and Eastern were associated with high competence words such as INTRICATE.

D.1 Warmth

Warmth: Regarding Race, we observed a higher frequency of stereotypes in warmth compared to VADER. Moreover, stereotype levels differed significantly across race. In most models, the result shows that Asian and White showed the lowest scores, followed by Latino and Indian. Black and Middle Eastern received the highest scores. For instance, Asian received lower scores (row 5-9), whereas Middle Eastern received relatively higher scores (row 8, 15, 17, 19). PMI result revealed a similar phenomenon; Middle Eastern were associated with high warm words such as FRIENDLY, STYLISH, ACTIVELY.

			Llama 3.2		Instruct BLIP		Pixtral		Qwen 2.5	
			11B	90B	7B	13B	12B	Large	7B	72B
Color	Blue vs. Red		0.21	1.96	7.19***	0.47	1.22	2.05*	3.03**	3.89***
	Blue vs. White		-1.40	0.10	6.45***	-1.11	3.38***	2.73**	0.70	-2.19*
	Red vs. White		-1.61	-1.86	-0.89	-1.59	2.18*	0.65	-2.31*	-6.07***
Gender	Male vs. Female		-6.41***	-3.66***	-8.63***	-8.29***	-15.26***	-7.89***	-2.35*	-3.30***
Race	Asian vs. Black		-2.20*	-0.46	-5.08***	-4.31***	-2.52*	-4.26***	3.21***	0.98
	Asian vs. Indian		0.91	2.28*	0.87	1.12	6.22***	2.72*	1.82	3.57***
	Asian vs. Latino		0.77	0.54	-3.32***	-1.60	0.01	-2.33*	1.19	0.68
	Asian vs. M.E.		0.27	1.49	0.95	-0.12	0.69	-0.83	2.33*	1.07
	Asian vs. White		0.55	-3.72***	-0.57	-1.32	2.36*	-3.45***	0.05	0.73
	Black vs. Indian		3.16***	2.67*	6.11***	5.35***	8.66***	6.92***	-1.38	2.52*
	Black vs. Latino		3.00***	1.01	1.87	2.79*	2.54*	1.83	-2.01	-0.30
	Black vs. M.E.		2.54*	1.90	6.05***	4.24***	3.19***	3.44***	-0.87	0.07
	Black vs. White		2.74*	-3.31***	4.59***	2.98***	4.79***	0.79	-3.16**	-0.26
	Indian vs. Latino		-0.14	-1.74	-4.20***	-2.68*	-6.21***	-4.99***	-0.61	-2.82***
	Indian vs. M.E.		-0.65	-0.79	0.10	-1.24	-5.40***	-3.56***	0.51	-2.47*
	Indian vs. White		-0.36	-5.93***	-1.44	-2.39*	-3.65***	-6.11***	-1.76	-2.75*
	Latino vs. M.E.		-0.50	0.96	4.25***	1.49	0.67	1.59	1.10	0.38
	Latino vs. White		-0.22	-4.32***	2.72*	0.28	2.34*	-1.05	-1.14	0.04
	M.E. vs. White		0.29	-5.32***	-1.52	-1.19	1.67	-2.60*	-2.25*	-0.33
Total			6/19	8/19	11/19	8/19	14/19	13/19	7/19	9/19

Table 7: Full experimental result for VADER metric

			Llama 3.2		Instruct BLIP		Pixtral		Qwen 2.5	
			11B	90B	7B	13B	12B	Large	7B	72B
Color	Blue vs. Red		-11.18***	-9.08**	-3.67***	-7.95**	-15.1***	-10.39**	-10.96***	-12.87**
	Blue vs. White		-11.08***	-9.43**	0.67	-6.39**	-21.36***	-14.52**	-11.33***	-16.16**
	Red vs. White		0.03	-0.21	4.45***	1.62	-5.66***	-3.78**	-0.40	-3.15**
Gender	Male vs. Female		-35.93***	-20.55**	-35.85***	-20.43**	-39.33***	-28.01**	-1.18	-4.51**
Race	Asian vs. Black		-0.82	3.57**	1.60	3.88**	4.76**	2.11*	0.94	3.72**
	Asian vs. Indian		4.19**	7.07**	4.07**	3.92**	13.35***	9.23**	0.37	7.10**
	Asian vs. Latino		0.01	4.39**	0.19	3.16**	9.19**	3.70**	0.66	3.38**
	Asian vs. M.E.		3.15**	5.08**	6.66**	5.46**	16.06***	11.23**	5.60**	10.54**
	Asian vs. White		-1.99*	-0.26	-0.50	1.09	8.64**	3.45**	-0.32	1.84
	Black vs. Indian		5.03***	3.62**	2.48*	0.05	8.63**	7.13**	-0.54	3.48**
	Black vs. Latino		0.82	0.88	-1.40	-0.71	4.58**	1.63	-0.28	-0.29
	Black vs. M.E.		3.85**	1.69	5.13**	1.72	11.72**	9.02**	4.73**	7.08**
	Black vs. White		-1.12	-3.90**	-2.08*	-2.78*	4.01**	1.39	-1.26	-1.79
	Indian vs. Latino		-4.10**	-2.74**	-3.79**	-0.75	-4.31**	-5.50**	0.27	-3.73**
	Indian vs. M.E.		-1.07	-1.90	2.57*	1.64	3.21**	1.88	5.14**	3.66**
	Indian vs. White		-5.99***	-7.31**	-4.47**	-2.78*	-4.81**	-5.90**	-0.68	-5.10**
	Latino vs. M.E.		3.04**	0.82	6.46**	2.40*	7.41**	7.51**	4.85**	7.19**
	Latino vs. White		-1.96	-4.66**	-0.69	-2.09*	-0.46	-0.25	-0.95	-1.48
	M.E. vs. White		-5.13**	-5.42**	-7.22**	-4.44**	-7.79**	-7.68**	-5.95**	-8.40**
Total			11/19	13/19	13/19	12/19	18/19	15/19	7/19	15/19

Table 8: Full experimental result for Competence metric in SCM

E Detailed Result on PMI

Tables from 10 to 31 show the detailed results of PMI analysis.

			Llama 3.2		Instruct BLIP		Pixtral		Qwen 2.5	
			11B	90B	7B	13B	12B	Large	7B	72B
Color	Blue	vs. Red	9.85***	11.96***	4.13***	9.80***	13.83***	8.80***	8.31***	8.40***
	Blue	vs. White	4.87***	9.71***	-2.36*	5.96***	0.47	0.24	-0.02	-5.85***
	Red	vs. White	-4.96***	-2.67**	-6.62***	-4.02***	-14.18***	-8.77***	-8.34***	-14.29***
Gender	Male	vs. Female	28.65***	23.72***	27.08***	35.62***	17.79***	12.09***	-6.05***	0.71
Race	Asian	vs. Black	-5.44***	-2.15*	-9.19***	-4.99***	-4.02***	-2.80**	0.91	-0.60
	Asian	vs. Indian	-4.79***	-5.25***	-4.87***	-1.40	3.33***	2.13*	2.07*	3.31***
	Asian	vs. Latino	-3.63***	-0.40	-5.02***	-3.14**	-3.20**	-0.04	1.11	-0.23
	Asian	vs. M.E.	-4.61***	-4.23***	-7.83***	-4.30***	-5.45***	-3.45***	-2.30*	-5.20***
	Asian	vs. White	-2.80**	2.40*	-3.11**	-1.31	-0.89	0.15	2.65**	2.51*
	Black	vs. Indian	0.72	-3.31***	4.35***	3.77***	7.32***	4.91***	1.11	3.87***
	Black	vs. Latino	1.86	1.75	4.24***	1.97*	0.77	2.77**	0.19	0.37
	Black	vs. M.E.	0.82	-2.17*	1.60	0.73	-1.37	-0.57	-3.23**	-4.52***
	Black	vs. White	2.59**	4.60***	6.07***	3.73***	3.09**	2.94**	1.75	3.19**
	Indian	vs. Latino	1.11	4.85***	-0.15	-1.74	-6.49***	9.23***	-0.91	-3.49***
	Indian	vs. M.E.	0.13	1.10	-2.86**	-3.01**	-8.88***	-5.56***	-4.32***	-8.51***
	Indian	vs. White	1.87	7.65***	1.69*	0.04	-4.30***	-1.97*	0.64	-0.61
	Latino	vs. M.E.	-0.97	-3.74***	-2.67**	-1.21	-2.17*	-3.32***	-3.34***	-4.95***
	Latino	vs. White	0.76	2.79**	1.83	1.82	2.29*	0.18	1.56	2.85**
	M.E.	vs. White	1.77	6.56***	4.60***	2.97**	4.47***	3.48***	4.95***	7.58***
Total			10/19	16/19	15/19	12/19	15/19	14/19	10/19	14/19

Table 9: Full experimental result for Competence metric in SCM

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
courage: 1.359 improvement: 1.292 intricately: 0.903 substantial: 0.804 intelligence: 0.774 satisfaction: 0.688 successful: 0.581 proficiency: 0.581 engaging: 0.536 achieve: 0.534 talents: 0.531 career: 0.494 accurate: 0.456 accomplished: 0.411 competitive: 0.344 words: 0.310 trust: 0.287 efficiently: 0.282 preparedness: 0.278 effective: 0.276	accuracy: 1.087 suitable: 1.087 capability: 0.987 straightforward: 0.920 productive: 0.894 words: 0.749 successful: 0.572 learning: 0.435 engaging: 0.380 efficiency: 0.350 powerful: 0.309 able: 0.263 dedicated: 0.250 job: 0.250 knowledgeable: 0.234 beauty: 0.233 role: 0.222 something: 0.216 trust: 0.210 responsibility: 0.203	rugged: 1.587 trust: 1.172 satisfaction: 0.780 character: 0.780 responsibility: 0.587 accurate: 0.558 satisfied: 0.555 performs: 0.537 hardworking: 0.517 accomplished: 0.489 diligently: 0.450 chosen: 0.407 skillful: 0.395 intently: 0.365 capabilities: 0.365 competence: 0.352 achieving: 0.324 accuracy: 0.320 secure: 0.312 enjoyable: 0.307	accurate: 0.988 able: 0.964 enthusiasm: 0.698 secure: 0.469 skilled: 0.423 engaging: 0.364 strength: 0.363 stable: 0.362 creativity: 0.353 assistant: 0.353 confidence: 0.332 precise: 0.189 equipped: 0.154 confidently: 0.143 confident: 0.128 efficiently: 0.124 ensuring: 0.108 functioning: 0.059 graceful: 0.054 role: 0.031
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
job: 0.673 words: 0.667 efficient: 0.580 dedicated: 0.512 determination: 0.464 flair: 0.442 knowledge: 0.395 beauty: 0.378 intricate: 0.359 enthusiastic: 0.357 distinguished: 0.357 productive: 0.300 secured: 0.287 skilled: 0.267 improvement: 0.258 suitable: 0.201 graceful: 0.165 role: 0.139 determined: 0.126 prepared: 0.125	accurately: 0.903 character: 0.785 substantial: 0.688 composure: 0.581 sleekly: 0.581 durable: 0.581 thorough: 0.503 easily: 0.466 proper: 0.450 efficient: 0.444 effectively: 0.333 job: 0.324 suitable: 0.287 enhance: 0.259 intricate: 0.254 powerful: 0.233 strength: 0.226 satisfaction: 0.208 efficiency: 0.203 trained: 0.189	job: 0.897 assistant: 0.874 confidence: 0.575 intricate: 0.301 disciplined: 0.295 confident: 0.231 effort: 0.188 thoughtful: 0.173 fully: 0.170 characterized: 0.160 strength: 0.142 skill: 0.051 equipped: 0.043 advanced: 0.039 something: 0.013 sophisticated: -0.010 intently: -0.010 task: -0.023 quality: -0.029 precise: -0.041	flair: 0.897 distinguished: 0.597 strong: 0.498 intricate: 0.431 job: 0.375 profession: 0.275 skill: 0.258 fully: 0.253 something: 0.124 elegant: 0.122 secured: 0.121 confidently: 0.098 learning: 0.095 skilled: 0.061 poised: 0.061 highly: 0.057 confident: 0.045 equipped: 0.033 strength: 0.032 engaging: 0.030

Table 10: Top 20 competence-related words associated with ‘Blue’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
charismatic: 1.329 hardworking: 0.870 flair: 0.733 accomplished: 0.722 stable: 0.663 functioning: 0.663 capabilities: 0.615 personality: 0.570 success: 0.570 distinguished: 0.570 talents: 0.520 commanding: 0.516 bright: 0.479 thoroughly: 0.471 prowess: 0.455 artistic: 0.435 ensuring: 0.427 committed: 0.400 characterized: 0.395 consistent: 0.374	character: 1.046 abilities: 1.010 assistant: 0.894 quality: 0.854 effective: 0.783 charismatic: 0.783 bright: 0.778 flair: 0.757 successful: 0.631 personality: 0.553 achieve: 0.494 precise: 0.471 preparedness: 0.442 career: 0.439 powerful: 0.368 performance: 0.345 competitive: 0.342 equipped: 0.319 capable: 0.299 job: 0.257	easily: 0.679 bright: 0.659 respected: 0.585 capabilities: 0.585 achieving: 0.585 enthusiastic: 0.494 aware: 0.470 advanced: 0.439 technological: 0.363 reliable: 0.363 eager: 0.343 fully: 0.342 highly: 0.322 automatic: 0.286 poised: 0.239 prepared: 0.232 talented: 0.207 abilities: 0.198 efficient: 0.188 ability: 0.175	abilities: 0.727 performs: 0.682 artistic: 0.682 beauty: 0.668 strong: 0.484 personality: 0.438 pleasing: 0.418 skilled: 0.383 successful: 0.353 suitable: 0.322 prepared: 0.315 accuracy: 0.253 talent: 0.253 precise: 0.236 trained: 0.213 powerful: 0.188 responsible: 0.181 easily: 0.170 profession: 0.166 determination: 0.128
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
competitive: 1.214 proper: 0.899 thoughtful: 0.860 satisfied: 0.729 powerful: 0.702 bright: 0.654 character: 0.577 strength: 0.539 intricately: 0.526 flair: 0.439 improvement: 0.392 productive: 0.297 secured: 0.284 equipped: 0.260 effort: 0.198 profession: 0.181 preparedness: 0.178 fully: 0.150 suitable: 0.133 elegant: 0.130	flair: 0.965 bright: 0.938 success: 0.813 competitive: 0.791 creativity: 0.720 diligently: 0.702 satisfied: 0.676 knowledge: 0.676 vivid: 0.633 stable: 0.550 highly: 0.527 strong: 0.516 improvement: 0.443 extensive: 0.357 specializing: 0.357 effort: 0.351 determined: 0.318 effectively: 0.302 attitude: 0.280 dedicated: 0.257	straightforward: 1.090 preparedness: 0.875 words: 0.839 learning: 0.576 strong: 0.576 highly: 0.558 difficult: 0.498 bright: 0.486 technological: 0.438 sophisticated: 0.387 characterized: 0.383 artistic: 0.351 speed: 0.313 elegant: 0.262 confidently: 0.254 secured: 0.237 suitable: 0.235 expertise: 0.213 strength: 0.213 powerful: 0.207	advanced: 1.160 words: 1.160 secure: 1.090 bright: 0.724 triumphant: 0.713 proficiency: 0.682 competitive: 0.575 strong: 0.476 fully: 0.438 effort: 0.327 secured: 0.294 performance: 0.292 elegant: 0.195 poised: 0.177 confidently: 0.172 equipped: 0.142 engaging: 0.111 confident: 0.106 skill: 0.103 learning: 0.073

Table 11: Top 20 competence-related words associated with ‘Red’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
improve: 1.018 skillfully: 0.925 achieved: 0.826 highly: 0.796 attractive: 0.796 meticulously: 0.710 talented: 0.626 delectable: 0.603 triumphant: 0.603 qualities: 0.603 secure: 0.545 knowledgeable: 0.539 enhance: 0.525 clean: 0.462 versatility: 0.419 powerful: 0.409 trusted: 0.381 versatile: 0.381 cleaner: 0.374 accuracy: 0.371	durable: 1.012 technological: 0.645 sophisticated: 0.552 difficult: 0.453 progress: 0.437 responsibility: 0.299 highly: 0.272 expert: 0.217 clean: 0.180 trained: 0.137 intently: 0.135 quiet: 0.131 graceful: 0.119 profession: 0.112 effort: 0.112 cleaner: 0.108 fully: 0.108 diligently: 0.108 committed: 0.093 competitive: 0.093	quiet: 1.167 thoroughly: 1.035 stable: 0.743 insights: 0.708 versatile: 0.695 consistent: 0.649 technological: 0.582 formidable: 0.582 talents: 0.540 competitive: 0.529 sophisticated: 0.463 efficiency: 0.439 effective: 0.412 skillful: 0.389 elegant: 0.364 achieve: 0.360 proficiency: 0.360 improving: 0.334 versatility: 0.328 enhance: 0.305	knowledge: 1.150 bright: 0.980 expertise: 0.768 determination: 0.658 trained: 0.565 effectively: 0.520 thoroughly: 0.465 versatile: 0.449 proper: 0.449 task: 0.337 poised: 0.328 creativity: 0.302 elegant: 0.263 skill: 0.218 profession: 0.205 artistic: 0.186 progress: 0.172 clean: 0.171 cleaner: 0.148 skills: 0.143
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
highly: 1.068 efficiency: 0.599 thorough: 0.599 clean: 0.555 meticulously: 0.532 graceful: 0.525 technological: 0.521 sophisticated: 0.498 disciplined: 0.461 skills: 0.376 strong: 0.376 prowess: 0.376 accuracy: 0.351 precise: 0.351 characterized: 0.316 expertise: 0.297 advanced: 0.294 poised: 0.290 quiet: 0.277 secure: 0.277	beauty: 0.777 knowledgeable: 0.732 meticulously: 0.655 clean: 0.632 intelligence: 0.625 disciplined: 0.568 technological: 0.552 proficiency: 0.525 competence: 0.515 capability: 0.509 sophisticated: 0.484 expertise: 0.474 enhanced: 0.402 intellectual: 0.394 quality: 0.385 characterized: 0.375 accurate: 0.362 intricately: 0.318 skill: 0.317 learning: 0.303	character: 1.342 triumphant: 0.827 clean : 0.638 expertise: 0.605 prepared: 0.511 consistent: 0.424 performance: 0.416 confidence: 0.342 skilled: 0.324 precise: 0.320 technological: 0.261 role: 0.250 professional: 0.197 intently: 0.172 cleaner: 0.165 task: 0.150 something: 0.134 graceful: 0.127 profession: 0.109 effort: 0.102	consistent: 1.583 disciplined: 0.955 intellectual: 0.824 clean: 0.726 skills: 0.665 technological: 0.641 durable: 0.639 suitable: 0.583 expertise: 0.583 artistic: 0.476 precise: 0.476 intently: 0.402 cleaner: 0.349 confidence: 0.329 sophisticated: 0.311 enthusiasm: 0.302 role: 0.297 skilled: 0.269 powerful: 0.261 quality: 0.184

Table 12: Top 20 competence-related words associated with ‘White (color)’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
trustworthy: 0.990 charismatic: 0.875 distinguished: 0.838 achieved: 0.575 triumphant: 0.575 talented: 0.567 thorough: 0.543 admiration: 0.543 respected: 0.538 straightforward: 0.521 courage: 0.505 characterized: 0.446 intellectual: 0.405 effort: 0.369 trust: 0.354 achieving: 0.353 responsibility: 0.349 enhance: 0.327 credibility: 0.312 experts: 0.312	qualities: 0.987 effort: 0.827 characterized: 0.681 effectively: 0.625 charismatic: 0.625 progress: 0.535 secured: 0.511 meticulously: 0.502 suitable: 0.502 responsibility: 0.492 enthusiasm: 0.402 intellectual: 0.376 words: 0.371 responsible: 0.347 straightforward: 0.335 achieve: 0.335 flair: 0.335 difficult: 0.328 equipped: 0.325 performance: 0.298	distinguished: 0.985 rugged: 0.985 speed: 0.763 preparedness: 0.763 intellectual: 0.705 hardworking: 0.680 smoothly: 0.570 secure: 0.549 bright: 0.545 responsibility: 0.533 accurately: 0.531 insights: 0.526 character: 0.500 capabilities: 0.500 satisfied: 0.400 eager: 0.390 accurate: 0.364 improve: 0.340 effective: 0.318 progress: 0.217	knowledge: 0.814 distinguished: 0.784 improve: 0.784 bright: 0.644 quality: 0.644 enhance: 0.592 character: 0.592 accurately: 0.592 personality: 0.547 easily: 0.436 accurate: 0.379 secure: 0.374 versatile: 0.306 proper: 0.306 trained: 0.296 thoroughly: 0.270 creativity: 0.229 something: 0.184 performance: 0.171 precise: 0.165
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
distinguished: 0.774 responsibility: 0.774 successful: 0.734 job: 0.675 satisfied: 0.634 skillfully: 0.511 productive: 0.494 thorough: 0.482 secured: 0.470 proper: 0.456 accuracy: 0.449 thoughtful: 0.380 intently: 0.342 intellectual: 0.292 enthusiasm: 0.291 precise: 0.278 words: 0.245 disciplined: 0.238 characterized: 0.233 triumphant: 0.219	distinguished: 0.999 speed: 0.736 performs: 0.736 accurately: 0.584 responsibility: 0.513 strongly: 0.513 easily: 0.468 secure: 0.439 progress: 0.342 character: 0.335 secured: 0.333 meticulously: 0.302 learning: 0.301 thoughtfulness: 0.298 disciplined: 0.253 enthusiastic: 0.238 enthusiasm: 0.227 flair: 0.221 intelligence: 0.221 performance: 0.216	extensive: 0.994 stable: 0.994 character: 0.994 intently: 0.487 learning: 0.479 powerful: 0.432 thoughtful: 0.428 triumphant: 0.409 skilled: 0.366 speed: 0.348 skill: 0.278 fully: 0.257 words: 0.257 disciplined: 0.229 confidence: 0.216 cleaner: 0.188 performance: 0.185 highly: 0.146 picking: 0.146 something: 0.143	enthusiastic: 0.797 secure: 0.767 triumphant: 0.575 distinguished: 0.575 intellectual: 0.437 intently: 0.426 learning: 0.362 competitive: 0.352 advanced: 0.312 skilled: 0.275 strong: 0.253 secured: 0.238 highly: 0.238 enthusiasm: 0.224 engaging: 0.222 skill: 0.207 effort: 0.201 something: 0.197 fully: 0.193 job: 0.182

Table 13: Top 20 competence-related words associated with ‘man’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
graceful: 1.010 improvement: 0.872 capable: 0.848 automatic: 0.817 capability: 0.762 strong: 0.659 secured: 0.657 handworking: 0.631 competence: 0.628 skillfully: 0.595 elegant: 0.568 competent: 0.565 accomplished: 0.540 accurately: 0.525 confidently: 0.523 determination: 0.505 satisfaction: 0.479 powerful: 0.453 determined: 0.449 reliable: 0.446	capability: 1.013 accomplished: 1.013 career: 0.920 capable: 0.903 commanding: 0.691 competent: 0.670 strong: 0.616 sophisticated: 0.598 indiscernible: 0.598 elegant: 0.582 bright: 0.563 competence: 0.550 intelligence: 0.527 beauty: 0.506 graceful: 0.506 assistant: 0.498 strength: 0.493 durable: 0.472 determination: 0.410 confident: 0.403	flair: 1.015 formidable: 0.948 competent: 0.933 proficiency: 0.908 graceful: 0.894 highly: 0.889 accomplished: 0.866 poised: 0.861 skillfully: 0.822 skillful: 0.822 capable: 0.819 whether: 0.808 suitable: 0.720 beauty: 0.717 assistant: 0.636 talented: 0.636 achieving: 0.600 trust: 0.600 competence: 0.580 strong: 0.567	assistant: 0.993 beauty: 0.894 poised: 0.807 dedicated: 0.771 elegant: 0.767 graceful: 0.740 expertise: 0.649 talent: 0.610 sophisticated: 0.508 determined: 0.501 enthusiasm: 0.491 capabilities: 0.479 role: 0.446 confidence: 0.395 able: 0.341 skilled: 0.315 artistic: 0.293 performs: 0.293 learning: 0.256 professional: 0.252
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
satisfaction: 1.003 enjoyable: 1.003 beauty: 0.939 efficient: 0.851 elegant: 0.813 graceful: 0.810 creativity: 0.781 intelligence: 0.740 competence: 0.716 confidence: 0.713 determination: 0.625 strong: 0.518 poised: 0.499 efficiency: 0.463 character: 0.418 diligently: 0.418 success: 0.418 confident: 0.367 improvement: 0.359 powerful: 0.351	assistant: 0.864 sleekly: 0.849 beauty: 0.785 graceful: 0.724 competence: 0.710 elegant: 0.682 efficient: 0.542 smoothly: 0.516 competent: 0.471 effectively: 0.454 skills: 0.454 poised: 0.434 strong: 0.433 suited: 0.416 composure: 0.416 confidence: 0.394 vivid: 0.373 determined: 0.365 determination: 0.364 satisfied: 0.349	secure: 1.006 elegant: 0.786 graceful: 0.475 assistant: 0.475 advanced: 0.470 poised: 0.458 technological: 0.421 sophisticated: 0.421 suitable: 0.402 consistent: 0.311 confidently: 0.301 strong: 0.269 strength: 0.251 artistic: 0.232 confident: 0.211 clean: 0.188 difficult: 0.150 role: 0.145 intricate: 0.136 quality: 0.122	poised: 0.853 elegant: 0.829 beauty: 0.788 disciplined: 0.623 confidence: 0.516 focused: 0.425 skills: 0.382 assistant: 0.332 words: 0.332 confidently: 0.325 flair: 0.310 proficiency: 0.310 intricately: 0.273 confident: 0.258 intricate: 0.233 determined: 0.161 artistic: 0.139 thoughtful: 0.136 suitable: 0.131 technological: 0.126

Table 14: Top 20 competence-related words associated with ‘woman’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
admiration: 0.929 efficiently: 0.684 preparedness: 0.669 respected: 0.621 trained: 0.488 consistent: 0.433 stable: 0.414 ensuring: 0.371 committed: 0.344 efficient: 0.324 capabilities: 0.321 difficult: 0.241 reliable: 0.234 responsibility: 0.228 prepared: 0.211 profession: 0.199 easily: 0.192 diligently: 0.192 dedicated: 0.183 equipped: 0.170	graceful: 1.079 precise: 1.072 consistent: 1.027 progress: 0.982 efficiency: 0.927 competitive: 0.638 characterized: 0.543 skill: 0.512 difficult: 0.413 artistic: 0.381 quiet: 0.371 intellectual: 0.316 job: 0.290 beauty: 0.260 quality: 0.249 role: 0.209 profession: 0.171 powerful: 0.150 responsibility: 0.143 capable: 0.125	satisfaction: 1.103 proficiency: 1.103 performs: 1.052 skillful: 0.910 intellectual: 0.823 advanced: 0.634 competence: 0.568 chosen: 0.501 progress: 0.464 graceful: 0.441 improve: 0.358 suitable: 0.350 respected: 0.347 responsible: 0.307 efficient: 0.293 equipped: 0.288 success: 0.281 hardworking: 0.255 secure: 0.234 skill: 0.232	proper: 0.893 functioning: 0.857 progress: 0.786 trained: 0.745 role: 0.667 effectively: 0.549 pleasing: 0.543 easily: 0.371 skill: 0.366 graceful: 0.340 profession: 0.327 powerful: 0.313 ensuring: 0.309 task: 0.274 determination: 0.272 efficiently: 0.272 prepared: 0.173 responsible: 0.171 precise: 0.167 something: 0.164
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
proficiency: 2.349 highly: 1.193 disciplined: 1.112 quiet: 0.665 accuracy: 0.646 enthusiasm: 0.618 characterized: 0.552 quality: 0.467 preparedness: 0.377 enthusiastic: 0.349 precise: 0.324 determined: 0.295 intellectual: 0.282 performance: 0.255 knowledge: 0.250 ensuring: 0.173 role: 0.162 suitable: 0.161 equipped: 0.135 prepared: 0.117	competent: 1.478 proper: 0.877 accuracy: 0.646 beauty: 0.594 durable: 0.594 satisfaction: 0.568 highly: 0.523 consistent: 0.424 words: 0.371 enthusiastic: 0.363 meticulously: 0.361 disciplined: 0.353 effort: 0.340 rugged: 0.308 precise: 0.304 preparedness: 0.297 prepared: 0.272 suitable: 0.210 skills: 0.153 dedicated: 0.149	job: 1.899 suitable: 0.651 graceful: 0.614 technological: 0.533 quality: 0.471 precise: 0.400 difficult: 0.329 secured: 0.321 bright: 0.312 advanced: 0.304 performance: 0.275 clean: 0.165 artistic: 0.159 elegant: 0.141 professional: 0.076 role: 0.064 engaging: 0.063 confidently: 0.049 something: 0.016 profession: 0.009	flair: 1.155 durable: 1.059 disciplined: 0.768 quality: 0.483 poised: 0.483 powerful: 0.434 intellectual: 0.396 suitable: 0.332 prepared: 0.306 performance: 0.271 role: 0.193 bright: 0.192 professional: 0.153 precise: 0.118 highly: 0.118 task: 0.115 skill: 0.107 clean: 0.113 expertise: 0.074 engaging: 0.019

Table 15: Top 20 competence-related words associated with ‘Asian’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
improvement: 2.181 success: 1.640 highly: 1.155 commanding: 1.055 effective: 0.865 powerful: 0.812 efficient: 0.699 secure: 0.697 accuracy: 0.671 successful: 0.640 strength: 0.618 versatility: 0.582 talented: 0.570 reliable: 0.553 capable: 0.545 competence: 0.487 artistic: 0.460 personality: 0.432 prepared: 0.408 confident: 0.406	effort: 0.928 skills: 0.521 strength: 0.504 competence: 0.378 confidence: 0.333 expert: 0.323 capable: 0.311 diligently: 0.298 confident: 0.297 characterized: 0.206 determination: 0.194 poised: 0.172 graceful: 0.158 personality: 0.151 enthusiasm: 0.113 professional: 0.076 confidently: 0.068 job: 0.047 role: 0.034 elegant: 0.022	reliable: 1.370 technologically: 1.214 advanced: 0.960 personality: 0.871 attitude: 0.802 accomplished: 0.791 enthusiastic: 0.678 chosen: 0.675 committed: 0.645 ability: 0.635 talents: 0.633 eager: 0.613 efficiency: 0.565 satisfied: 0.560 prepared: 0.528 hardworking: 0.522 achieve: 0.522 determined: 0.513 strong: 0.463 effectively: 0.460	pleasing: 1.157 skilled: 0.877 determined: 0.648 talent: 0.648 stable: 0.601 role: 0.530 successful: 0.485 confidence: 0.448 engaging: 0.371 confidently: 0.332 poised: 0.318 strength: 0.303 powerful: 0.275 strong: 0.233 efficiently: 0.233 confident: 0.196 proper: 0.177 precise: 0.129 professional: 0.117 easily: 0.101
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
improvement: 1.272 prowess: 1.109 skills: 1.109 productive: 1.092 elegant: 0.926 advanced: 0.787 strength: 0.768 dedicated: 0.659 artistic: 0.573 engaging: 0.542 determination: 0.479 discipline: 0.457 confident: 0.432 triumphant: 0.424 thoughtful: 0.393 learning: 0.311 effort: 0.309 preparedness: 0.281 confidently: 0.230 sophisticated: 0.223	sleekly: 1.603 improvement: 1.255 flair: 1.018 strong: 0.878 powerful: 0.603 elegant: 0.585 enhanced: 0.533 quality: 0.516 poised: 0.451 attitude: 0.416 progress: 0.416 strength: 0.408 enthusiasm: 0.366 confidence: 0.347 confidently: 0.339 sophisticated: 0.336 highly: 0.327 bright: 0.300 prepared: 0.298 enhance: 0.281	intently: 0.851 thoughtful: 0.718 effort: 0.688 sophisticated: 0.680 technological: 0.561 confident: 0.530 engaging: 0.454 artistic: 0.445 highly: 0.436 performance: 0.270 bright: 0.244 speed: 0.243 profession: 0.179 difficult: 0.165 role: 0.107 cleaner: 0.081 professional: 0.037 quality: 0.030 something: 0.029 strength: 0.021	strong: 1.033 enthusiasm: 0.852 confidently: 0.582 engaging: 0.466 skilled: 0.423 bright: 0.418 confident: 0.406 fully: 0.381 skill: 0.361 poised: 0.345 intently: 0.338 expertise: 0.328 confidence: 0.248 elegant: 0.238 cleaner: 0.214 highly: 0.203 something: 0.180 professional: 0.115 strength: 0.109 precise: 0.079

Table 16: Top 20 competence-related words associated with 'Black'

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
graceful: 1.675 thoroughly: 1.246 admiration: 0.983 thorough: 0.983 accurate: 0.791 determined: 0.773 enhance: 0.642 progress: 0.630 flair: 0.568 elaborate: 0.544 artistic: 0.433 straightforward: 0.305 personality: 0.257 elegant: 0.225 respected: 0.190 meticulously: 0.190 quiet: 0.147 competence: 0.142 competent: 0.140 words: 0.139	engaging: 1.480 intricate: 0.892 elegant: 0.721 abilities: 0.673 knowledge: 0.451 quiet: 0.444 learning: 0.436 role: 0.416 intellectual: 0.351 artistic: 0.319 powerful: 0.310 responsibility: 0.304 able: 0.301 confidently: 0.275 expert: 0.267 prepared: 0.240 something: 0.207 skills: 0.187 fully: 0.181 preparedness: 0.162	distinguished: 2.174 stable: 0.823 intricate: 0.782 intently: 0.656 accurate: 0.589 diligently: 0.545 beauty: 0.461 knowledgeable: 0.267 responsibility: 0.210 artistic: 0.198 fully: 0.190 effective: 0.156 learning: 0.151 enthusiastic: 0.144 determination: 0.132 strong: 0.123 career: 0.121 elegant: 0.102 graceful: 0.094 powerful: 0.091	expertise: 0.711 confidence: 0.479 poised: 0.349 strength: 0.333 prepared: 0.326 determination: 0.264 profession: 0.226 job: 0.139 confidently: 0.123 responsible: 0.105 performance: 0.104 suitable: 0.100 efficiently: 0.094 graceful: 0.091 powerful: 0.083 diligently: 0.052 something: 0.032 elegant: 0.022 equipped: -0.019 professional: -0.057
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
quiet: 1.269 graceful: 1.176 beauty: 1.067 secured: 0.783 effort: 0.765 intricately: 0.733 intricate: 0.691 thoughtful: 0.652 fully: 0.624 dedicated: 0.475 suitable: 0.364 skilled: 0.324 strength: 0.305 meticulously: 0.280 job: 0.269 performance: 0.248 sophisticated: 0.188 enthusiastic: 0.105 ensuring: 0.082 profession: 0.080	intricately: 0.934 intricate: 0.825 vivid: 0.816 graceful: 0.775 learning: 0.482 thoughtful: 0.412 words: 0.360 intellectual: 0.352 disciplined: 0.341 characterized: 0.292 skill: 0.250 intently: 0.199 profession: 0.184 something: 0.151 accuracy: 0.149 determined: 0.094 artistic: 0.054 competence: 0.032 meticulously: 0.027 confidently: 0.026	intricate: 1.046 skill: 1.043 expertise: 0.981 difficult: 0.903 highly: 0.686 graceful: 0.603 confidently: 0.585 precise: 0.389 suitable: 0.318 quality: 0.290 consistent: 0.286 effort: 0.286 thoughtful: 0.164 equipped: 0.123 fully: 0.074 something: 0.066 sophisticated: 0.055 cleaner: 0.041 artistic: 0.027 prepared: 0.026	learning: 1.082 intricate: 0.899 graceful: 0.824 intently: 0.819 poised: 0.534 quality: 0.311 expertise: 0.295 artistic: 0.238 confident: 0.154 skilled: 0.048 technological: -0.001 professional: -0.018 role: -0.025 something: -0.053 performance: -0.086 secured: -0.094 bright: -0.124 equipped: -0.133 skill: -0.142 task: -0.161

Table 17: Top 20 competence-related words associated with ‘Indian’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
character: 1.331 hardworking: 1.315 effort: 0.957 able: 0.870 powerful: 0.651 competitive: 0.649 talents: 0.643 versatility: 0.635 preparedness: 0.560 straightforward: 0.524 durable: 0.506 creativity: 0.497 abilities: 0.464 precise: 0.446 knowledge: 0.434 ensuring: 0.425 talent: 0.425 reliable: 0.413 responsibility: 0.408 efficiency: 0.399	assistant: 1.656 durable: 1.241 knowledgeable: 1.037 responsibility: 1.023 career: 0.978 words: 0.939 trained: 0.656 difficult: 0.556 progress: 0.540 competitive: 0.518 secured: 0.463 skills: 0.351 something: 0.239 enthusiasm: 0.178 cleaner: 0.178 equipped: 0.130 performance: 0.099 responsible: 0.098 task: 0.096 confident: 0.094	specializing: 1.174 accurately: 0.965 versatility: 0.805 diligently: 0.715 talented: 0.696 accuracy: 0.654 efficiency: 0.562 consistent: 0.519 enthusiasm: 0.512 productive: 0.460 powerful: 0.397 able: 0.388 respected: 0.348 sophisticated: 0.339 fully: 0.316 hardworking: 0.256 accomplished: 0.243 expert: 0.232 profession: 0.226 responsibility: 0.211	enthusiasm: 1.509 creativity: 1.333 precise: 0.755 suitable: 0.696 engaging: 0.537 secure: 0.449 talent: 0.274 skills: 0.174 clean: 0.172 strong: 0.137 job: 0.102 confident: 0.092 diligently: 0.063 professional: 0.037 something: 0.024 determination: 0.011 prepared: -0.005 easily: -0.019 task: -0.030 cleaner: -0.072
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
character: 2.002 flair: 1.449 job: 0.680 confidence: 0.587 words: 0.587 knowledge: 0.528 competence: 0.371 intently: 0.340 determination: 0.208 poised: 0.208 elegant: 0.199 advanced: 0.194 clean: 0.194 confident: 0.176 cleaner: 0.175 fully: 0.160 equipped: 0.149 role: 0.113 expertise: 0.109 effort: 0.108	thoughtfulness: 1.208 easily: 1.208 attitude: 1.036 flair: 1.001 strong: 0.967 competitive: 0.934 satisfaction: 0.660 skills: 0.660 enthusiastic: 0.655 character: 0.629 trained: 0.516 fully: 0.512 enthusiasm: 0.480 strength: 0.313 determination: 0.264 performance: 0.204 secured: 0.169 confidence: 0.140 task: 0.134 expertise: 0.107	poised: 0.941 engaging: 0.742 powerful: 0.650 secured: 0.500 prepared: 0.411 effort: 0.324 profession: 0.307 thoughtful: 0.202 performance: 0.198 confident: 0.181 strength: 0.171 intently: 0.171 elegant: 0.168 consistent: 0.102 task: 0.076 advanced: 0.068 professional: 0.050 role: 0.047 cleaner: 0.043 bright: 0.029	job: 1.369 flair: 1.213 fully: 0.547 quality: 0.541 suitable: 0.527 intellectual: 0.454 strength: 0.382 durable: 0.269 bright: 0.250 confidence: 0.222 equipped: 0.206 prepared: 0.190 technological: 0.176 clean: 0.150 sophisticated: 0.132 confident: 0.112 elegant: 0.096 intently: 0.089 engaging: 0.077 artistic: 0.076

Table 18: Top 20 competence-related words associated with ‘Latino’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
attractive: 1.098 highly: 1.098 characterized: 1.039 functioning: 0.999 triumphant: 0.906 advanced: 0.818 career: 0.818 flair: 0.715 intricate: 0.624 precise: 0.584 engaging: 0.584 versatility: 0.525 competent: 0.464 elegant: 0.457 productive: 0.455 sophisticated: 0.391 abilities: 0.354 difficult: 0.311 beauty: 0.298 proper: 0.298	able: 1.142 highly: 0.929 intricate: 0.913 bright: 0.794 quality: 0.695 dedicated: 0.681 quiet: 0.636 creativity: 0.579 performance: 0.574 preparedness: 0.562 knowledge: 0.472 personality: 0.447 fully: 0.386 beauty: 0.355 artistic: 0.341 determined: 0.316 characterized: 0.310 elegant: 0.306 profession: 0.305 expert: 0.289	technological: 1.095 character: 1.095 intricate: 0.966 graceful: 0.533 improve: 0.457 productive: 0.451 artistic: 0.407 beauty: 0.399 intently: 0.358 able: 0.320 performs: 0.308 automatic: 0.281 strength: 0.243 competent: 0.243 versatility: 0.212 functioning: 0.205 powerful: 0.166 precise: 0.144 engaging: 0.123 competitive: 0.113	artistic: 1.214 beauty: 1.008 accurate: 0.827 secure: 0.639 stable: 0.446 strong: 0.396 expertise: 0.356 poised: 0.356 graceful: 0.339 powerful: 0.313 role: 0.219 skill: 0.204 something: 0.196 elegant: 0.185 performance: 0.142 effectively: 0.133 task: 0.128 confidently: 0.047 determination: 0.008 diligently: -0.019
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
intricate: 0.870 secured: 0.679 words: 0.593 knowledge: 0.534 competence: 0.510 intellectual: 0.455 enthusiasm: 0.446 triumphant: 0.423 intricately: 0.320 profession: 0.219 confidence: 0.217 bright: 0.176 quality: 0.167 poised: 0.115 effort: 0.115 intently: 0.113 role: 0.093 engaging: 0.090 graceful: 0.078 preparedness: 0.057	character: 0.945 secured: 0.932 diligently: 0.749 intricate: 0.723 progress: 0.562 learning: 0.498 fully: 0.445 dedicated: 0.442 secure: 0.375 sophisticated: 0.374 intricately: 0.365 consistent: 0.354 rugged: 0.312 thoughtful: 0.276 enhance: 0.276 enthusiasm: 0.216 meticulous: 0.213 confidence: 0.191 ensuring: 0.179 skills: 0.157	characterized: 1.165 fully: 0.789 intricate: 0.676 speed: 0.580 elegant: 0.449 skill: 0.379 strength: 0.285 thoughtful: 0.178 confident: 0.142 confidently: 0.052 advanced: 0.044 role: 0.008 equipped: -0.018 professional: -0.027 something: -0.029 graceful: -0.120 clean: -0.137 technological: -0.142 consistent: -0.185 effort: -0.185	skills: 1.326 secured: 1.231 intricate: 0.844 disciplined: 0.841 confidently: 0.834 powerful: 0.507 confidence: 0.459 highly: 0.414 profession: 0.241 graceful: 0.140 strength: 0.097 artistic: 0.091 poised: 0.070 elegant: 0.070 quality: 0.070 equipped: 0.003 sophisticated: -0.023 clean: -0.026 role: -0.035 task: -0.050

Table 19: Top 20 competence-related words associated with ‘Middle Eastern’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
achieving: 1.296 triumphant: 1.104 rugged: 1.016 meticulously: 0.625 committed: 0.612 capability: 0.593 secure: 0.575 productive: 0.559 effort: 0.519 accurate: 0.519 successful: 0.519 ensuring: 0.491 able: 0.454 stable: 0.419 responsibility: 0.403 sophisticated: 0.393 equipped: 0.375 efficient: 0.326 capabilities: 0.326 assistant: 0.317	achieve: 1.268 abilities: 1.083 committed: 0.946 enthusiasm: 0.579 strong: 0.563 intellectual: 0.531 determined: 0.512 responsible: 0.484 effort: 0.479 secured: 0.405 diligently: 0.375 efficiency: 0.335 prepared: 0.328 personality: 0.327 competent: 0.318 bright: 0.298 progress: 0.290 job: 0.276 profession: 0.245 personality: 0.241	insights: 1.433 highly: 1.308 thoroughly: 1.130 effort: 1.020 bright: 0.908 responsibility: 0.870 versatile: 0.821 proficient: 0.805 consistent: 0.763 achieve: 0.638 accurate: 0.571 poised: 0.554 expert: 0.535 competitive: 0.443 talents: 0.441 equipped: 0.440 formidable: 0.433 suitable: 0.405 effective: 0.401 progress: 0.387	abilities: 1.739 talent: 0.680 equipped: 0.441 strong: 0.390 cleaner: 0.372 suitable: 0.364 responsible: 0.328 ensuring: 0.325 skill: 0.281 clean: 0.236 proper: 0.208 precise: 0.161 poised: 0.127 effectively: 0.127 diligently: 0.127 engaging: 0.113 determined: 0.095 graceful: 0.092 skills: 0.089 professional: 0.076
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
strong: 1.351 skills: 1.088 proper: 0.896 precise: 0.481 quality: 0.469 determination: 0.458 fully: 0.458 technological: 0.413 skilled: 0.394 prepared: 0.305 intellectual: 0.284 improvement: 0.252 something: 0.222 triumphant: 0.211 clean: 0.202 learning: 0.191 advanced: 0.181 determined: 0.181 bright: 0.179 expertise: 0.169	efficiency: 1.432 intelligence: 1.285 knowledgeable: 1.169 rugged: 0.963 vivid: 0.782 trained: 0.741 durable: 0.718 technological: 0.705 dedicated: 0.634 job: 0.634 thorough: 0.622 proper: 0.609 skilled: 0.460 suitable: 0.449 preparedness: 0.428 enhance: 0.419 accuracy: 0.378 cleaner: 0.336 secure: 0.326 quality: 0.308	graceful: 0.461 intently: 0.407 advanced: 0.304 consistent: 0.297 equipped: 0.228 task: 0.226 clean: 0.218 powerful: 0.207 prepared: 0.162 cleaner: 0.154 something: 0.092 sophisticated: 0.066 bright: 0.001 strength: -0.008 professional: -0.011 secured: -0.042 quality: -0.056 suitable: -0.086 elegant: -0.139 technological: -0.146	strong: 0.994 effort: 0.916 durable: 0.520 confidence: 0.432 sophisticated: 0.409 prepared: 0.387 suitable: 0.377 precise: 0.356 engaging: 0.344 skilled: 0.306 clean: 0.199 highly: 0.164 task: 0.131 elegant: 0.123 fully: 0.119 cleaner: 0.113 something: 0.044 equipped: 0.022 professional: 0.016 profession: 0.002

Table 20: Top 20 competence-related words associated with ‘White(race)’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
gratitude: 1.122 prayerful: 0.953 lively: 0.881 trusting: 0.733 appreciation: 0.688 active: 0.581 enjoys: 0.491 loving: 0.444 cheerful: 0.318 intimate: 0.301 welcoming: 0.292 trust: 0.287 admiration: 0.259 respect: 0.259 trustworthy: 0.259 enjoying: 0.222 respected: 0.203 authenticity: 0.203 attired: 0.200 gentle: 0.191	mutual: 1.157 suitable: 1.087 enjoys: 0.894 enjoying: 0.648 safe: 0.572 friendly: 0.454 warmth: 0.312 gentle: 0.281 calmness: 0.250 knowledgeable: 0.234 beauty: 0.233 trust: 0.210 respect: 0.194 lively: 0.194 confident: 0.179 secured: 0.157 warm: 0.119 committed: 0.113 kind: 0.093 stylish: 0.070	rugged: 1.587 trust: 1.172 active: 0.751 hot: 0.587 prayerful: 0.500 skillful: 0.395 cheerful: 0.350 secure: 0.312 enjoyable: 0.307 lively: 0.307 friendly: 0.242 enjoying: 0.234 strong: 0.231 quality: 0.227 smooth: 0.212 passionate: 0.211 suitable: 0.193 committed: 0.186 actively: 0.186 committed: 0.183	good: 1.130 natural: 0.915 secure: 0.469 welcoming: 0.214 comfortable: 0.201 enjoying: 0.193 confident: 0.128 taste: 0.080 delicious: 0.049 strong: 0.004 neutral: -0.020 prayerful: -0.062 manner: -0.064 appealing: -0.067 suitable: -0.285 kind: -0.435 fresh: -0.455 stylish: -0.574 pleasing: -0.657 actively: -0.279
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
solid: 1.442 heated: 1.094 enjoying: 0.855 comfort: 0.480 pleasant: 0.464 gentle: 0.454 grayish: 0.442 beauty: 0.378 enthusiastic: 0.357 prayerful: 0.317 manner: 0.288 secured: 0.287 lively: 0.280 welcoming: 0.269 respect: 0.258 friendly: 0.249 warm: 0.222 suitable: 0.201 cheerful: 0.183 attired: 0.120	heated: 0.733 respect: 0.581 friendliness: 0.581 enjoying: 0.329 suitable: 0.287 friendly: 0.250 welcoming: 0.203 thoughtful: 0.177 enthusiastic: 0.174 attitude: 0.130 pleasant: 0.113 warm: 0.110 secured: 0.106 lively: 0.096 cheerful: 0.068 confident: 0.064 actively: 0.047 good: 0.044 plain: 0.002 rugged: -0.004	hot: 1.089 respect: 1.014 solid: 0.767 active: 0.745 gentle: 0.382 confident: 0.231 neutral: 0.225 thoughtful: 0.173 kind: 0.169 actively: 0.141 stylish: 0.137 natural: 0.10789 smooth: -0.010 friendly: -0.025 quality: -0.029 bright: -0.060 plain: -0.073 enjoying: -0.077 fresh: -0.151 secured: -0.151	solid: 1.083 pleasant: 0.960 respect: 0.919 enjoying: 0.767 heated: 0.597 neutral: 0.545 strong: 0.498 comfortable: 0.396 kind: 0.375 welcoming: 0.346 natural: 0.265 gentle: 0.253 fine: 0.219 secured: 0.121 warm: 0.095 confident: 0.045 plain: 0.017 cheerful: -0.031 good: -0.031 stylish: -0.058

Table 21: Top 20 warmth-related words associated with ‘blue’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
flattering: 0.933 hot: 0.860 bright: 0.479 tasting: 0.433 committed: 0.400 welcoming: 0.359 mutual: 0.333 cheerful: 0.307 friendly: 0.305 beauty: 0.300 good: 0.285 secured: 0.240 fine: 0.213 fresh: 0.193 respected: 0.192 rugged: 0.183 natural: 0.160 strong: 0.149 quality: 0.135 confident: 0.126	quality: 0.854 fine: 0.854 intimate: 0.783 good: 0.783 bright: 0.716 fresh: 0.631 comfort: 0.631 respectful: 0.561 lively: 0.516 welcoming: 0.479 solid: 0.337 warmth: 0.302 natural: 0.233 comfortable: 0.194 confident: 0.186 strong: 0.179 friendliness: 0.125 manner: 0.125 passionate: 0.124 warm: 0.114	authenticity: 0.848 respectful: 0.848 pleasant: 0.778 bright: 0.659 appreciation: 0.642 respected: 0.585 loving: 0.585 hot: 0.585 enthusiastic: 0.494 reliable: 0.363 fresh: 0.332 comfort: 0.315 natural: 0.264 talented: 0.207 stylish: 0.199 committed: 0.152 safe: 0.136 knowledgeable: 0.094 prayerful: 0.083 comfortable: 0.066	safely: 0.727 comfort: 0.701 beauty: 0.668 fresh: 0.642 kind: 0.525 strong: 0.494 pleasing: 0.418 delicious: 0.353 safe: 0.353 actively: 0.324 suitable: 0.322 smooth: 0.206 plain: 0.090 stylish: 0.086 manner: -0.029 enjoying: -0.049 welcoming: -0.049 secure: -0.057 confident: -0.091 appealing: -0.107
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
thoughtful: 0.860 bright: 0.654 comfortable: 0.622 safe: 0.354 secured: 0.284 lively: 0.277 fine: 0.255 cheerful: 0.161 warm: 0.142 suitable: 0.133 fresh: 0.131 strong: 0.091 respectful: 0.074 good: 0.068 actively: 0.054 beauty: 0.053 manner: 0.047 confident: 0.024 respect: -0.008 secure: -0.008	bright: 0.938 safely: 0.830 affectionately: 0.702 solid: 0.633 hot: 0.587 warmth: 0.550 respectful: 0.550 strong: 0.516 appealing: 0.398 active: 0.398 comfortable: 0.311 attitude: 0.280 comfort: 0.213 secure: 0.202 good: 0.179 thoughtfulness: 0.172 cheerful: 0.162 fresh: 0.161 enjoying: 0.117 lively: 0.105	strong: 0.576 enjoying: 0.576 bright: 0.486 tasting: 0.406 good: 0.383 stylish: 0.361 comfort: 0.353 cheerful: 0.328 friendly: 0.254 secured: 0.237 suitable: 0.235 fresh: 0.213 fine: 0.197 manner: 0.181 thoughtfulness: 0.174 solid: 0.171 quality: 0.158 actively: 0.079 welcoming: 0.045 kind: -0.065	secure: 1.090 lightly: 0.954 warmth: 0.727 bright: 0.724 strong: 0.476 lively: 0.319 friendly: 0.302 cheerful: 0.295 warm: 0.295 secured: 0.294 tasting: 0.197 fresh: 0.193 good: 0.188 welcoming: 0.179 gentle: 0.116 confident: 0.106 kind: 0.090 actively: -0.021 manner: -0.030 stylish: -0.050

Table 22: Top 20 warmth-related words associated with ‘red’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
pure: 1.603 safely: 0.987 skillfully: 0.925 attractive: 0.796 safe: 0.755 talented: 0.626 qualities: 0.603 fluffy: 0.571 fine: 0.568 secure: 0.545 knowledgeable: 0.539 slicked: 0.483 tasting: 0.466 trusted: 0.381 rugged: 0.323 appealing: 0.298 calmness: 0.294 reliable: 0.216 plain: 0.199 respectful: 0.199	fluffy: 0.874 tasting: 0.645 fine: 0.289 smooth: 0.267 calmness: 0.230 delicious: 0.230 actively: 0.214 slicked: 0.174 committed: 0.093 neutral: 0.077 plain: 0.047 strong: 0.022 secured: 0.012 trust: -0.033 kind: -0.041 gentle: -0.074 attired: -0.091 stylish: -0.102 passionate: -0.115 manner: -0.124	harmonious: 0.997 taste: 0.582 skillful: 0.389 plain: 0.309 manner: 0.281 appealing: 0.252 smooth: 0.245 safely: 0.195 delicious: 0.180 welcoming: 0.145 neutral: 0.142 warm: 0.142 good: 0.057 friendly: 0.028 suitable: 0.013 quality: -0.010 respect: -0.032 strong: -0.045 lively: -0.046 attitude: -0.060	friendly: 1.005 bright: 0.980 comfortable: 0.565 prayerful: 0.565 appreciation: 0.427 taste: 0.407 plain: 0.342 stylish: 0.329 smooth: 0.311 neutral: 0.159 appealing: 0.156 manner: 0.087 safe: 0.079 pleasing: 0.029 confident: -0.042 suitable: -0.113 actively: -0.121 enjoying: -0.159 welcoming: -0.185 kind: -0.293
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
safe: 0.528 fine: 0.499 fluffy: 0.477 grayish: 0.461 fluffy: 0.429 strong: 0.376 smooth: 0.370 intimate: 0.335 secure: 0.277 quality: 0.260 lightly: 0.255 natural: 0.165 stylish: 0.157 neutral: 0.145 respectful: 0.096 confident: 0.075 plain: 0.037 actively: -0.029 fresh: -0.068 active: -0.069	calmness: 1.362 beauty: 0.777 knowledgeable: 0.732 slicked: 0.709 fluffy: 0.677 tasting: 0.677 fine: 0.437 quality: 0.385 gentle: 0.353 smooth: 0.331 prayerful: 0.303 thoughtfulness: 0.246 enthusiastic: 0.218 strong: 0.198 kind: 0.180 lightly: 0.148 rugged: 0.147 secure: 0.139 manner: 0.138 stylish: 0.131	fluffy: 0.868 smooth: 0.757 fine: 0.489 heated: 0.342 tasting: 0.242 welcoming: 0.226 good: 0.190 plain: 0.174 comfort: 0.119 natural: 0.050 manner: -0.019 neutral: -0.048 cheerful: -0.058 suitable: -0.058 confident: -0.069 active: -0.074 gentle: -0.074 fresh: -0.091 secured: -0.121 kind: -0.123	slicked: 1.583 fluffy: 0.980 smooth: 0.835 suitable: 0.583 nice: 0.467 active: 0.261 tasting: 0.204 manner: 0.187 quality: 0.184 stylish: 0.102 actively: 0.077 plain: 0.042 natural: -0.113 comfortable: -0.134 lively: -0.143 fresh: -0.153 confident: -0.164 good: -0.183 friendly: -0.188 neutral: -0.282

Table 23: Top 20 warmth-related words associated with ‘white (color)’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
trustworthy: 0.990 slicked: 0.990 gratitude: 0.853 harmonious: 0.768 coolness: 0.727 attired: 0.581 flavorful: 0.575 talented: 0.567 admiration: 0.543 respected: 0.538 tasting: 0.531 hot: 0.521 lively: 0.460 trustworthiness: 0.405 respect: 0.373 prayerful: 0.362 trust: 0.354 credibility: 0.312 authenticity: 0.290 welcoming: 0.286	qualities: 0.987 slicked: 0.872 welcoming: 0.665 attired: 0.590 secured: 0.511 suitable: 0.502 enjoying: 0.332 trust: 0.277 safe: 0.250 respect: 0.234 quality: 0.210 neutral: 0.190 kind: 0.176 stylish: 0.163 passionate: 0.139 kindness: 0.139 solid: 0.117 smooth: 0.094 lively: 0.094 gentle: 0.085	rugged: 0.985 flavorful: 0.815 solid: 0.763 smoothly: 0.570 secure: 0.549 bright: 0.545 lively: 0.513 friendly: 0.497 good: 0.442 neutral: 0.416 gratitude: 0.307 appreciation: 0.248 loving: 0.248 pleasantries: 0.248 taste: 0.248 plain: 0.215 stylish: 0.211 pleasant: 0.178 enthusiastic: 0.125 enjoying: 0.112	good: 0.784 friendly: 0.659 bright: 0.644 quality: 0.644 lively: 0.592 appreciation: 0.547 plain: 0.521 pleasant: 0.521 flavorful: 0.521 prayerful: 0.466 secure: 0.374 stylish: 0.309 taste: 0.236 delicious: 0.199 warm: 0.159 comfort: 0.132 comfortable: 0.122 strong: 0.083 safe: 0.007 neutral: -0.106
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
freshness: 0.774 good: 0.595 good: 0.582 fine: 0.549 respect: 0.549 skillfully: 0.511 taste: 0.511 secured: 0.470 prayerful: 0.412 lively: 0.389 thoughtful: 0.380 cheerful: 0.235 respectful: 0.231 enthusiastic: 0.190 quality: 0.156 tasting: 0.149 welcoming: 0.121 gentle: 0.112 activate: 0.100 secure: 0.090	friendliness: 0.806 slicked: 0.624 good: 0.479 secure: 0.439 solid: 0.371 lively: 0.338 kind: 0.338 secured: 0.333 respectful: 0.321 thoughtfulness: 0.298 welcoming: 0.249 enthusiastic: 0.238 friendly: 0.178 attitude: 0.165 fine: 0.165 hot: 0.106 safely: 0.098 thoughtful: 0.092 comfort: 0.058 prayerful: 0.055	lively: 0.731 prayerful: 0.579 respect: 0.464 thoughtful: 0.428 welcoming: 0.294 fine: 0.294 solid: 0.237 natural: 0.194 good: 0.164 secured: 0.091 kind: 0.073 plain: 0.025 fresh: 0.021 gentle: -0.006 active: -0.006 friendly: -0.035 bright: -0.096 actively: -0.112 neutral: -0.080 quality: -0.132	slicked: 0.990 enthusiastic: 0.797 secure: 0.767 prayerful: 0.767 respect: 0.575 lively: 0.518 warmth: 0.405 welcoming: 0.390 comfortable: 0.373 cheerful: 0.317 good: 0.295 tasting: 0.289 strong: 0.253 secured: 0.238 kind: 0.182 friendly: 0.176 quality: 0.132 actively: 0.113 natural: 0.112 smooth: 0.050

Table 24: Top 20 warmth-related words associated with ‘man’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
loving: 0.872 generous: 0.788 flattering: 0.788 strong: 0.659 secured: 0.657 skillfully: 0.595 reliable: 0.446 cheerful: 0.425 trusting: 0.425 kindness: 0.425 confident: 0.418 fluffy: 0.376 bright: 0.373 warm: 0.369 secure: 0.366 enjoys: 0.358 comfortable: 0.323 beauty: 0.289 passionate: 0.262 smooth: 0.230	intimate: 0.843 heated: 0.750 strong: 0.616 indiscernible: 0.598 tasting: 0.565 bright: 0.563 warm: 0.527 beauty: 0.506 enjoys: 0.472 confident: 0.403 calmness: 0.391 committed: 0.361 manner: 0.297 fresh: 0.261 natural: 0.242 fine: 0.235 friendliness: 0.221 comfort: 0.183 knowledgeable: 0.172 delicious: 0.106	skillfully: 0.822 skillful: 0.822 suitable: 0.720 beauty: 0.717 talented: 0.636 trust: 0.600 strong: 0.567 cheerful: 0.555 knowledgeable: 0.466 active: 0.458 harmonious: 0.430 attractive: 0.352 confident: 0.334 comfortable: 0.318 natural: 0.278 authenticity: 0.278 respectful: 0.278 delicious: 0.266 enjoyable: 0.249 welcoming: 0.240	beauty: 0.894 actively: 0.346 manner: 0.336 welcoming: 0.307 fresh: 0.301 natural: 0.293 appealing: 0.225 pleasing: 0.223 smooth: 0.209 suitable: 0.186 enjoying: 0.162 safely: 0.145 kind: 0.135 confident: 0.106 neutral: 0.098 safe: -0.007 strong: -0.088 comfortable: -0.132 comfort: -0.144 delicious: -0.229
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
enjoyable: 1.003 beauty: 0.939 heated: 0.781 warm: 0.645 comfort: 0.556 safe: 0.518 strong: 0.518 lightly: 0.482 manner: 0.411 comfortable: 0.371 confident: 0.367 grayish: 0.351 enjoying: 0.279 fluffy: 0.253 pleasant: 0.210 actively: 0.202 neutral: 0.155 natural: 0.146 solid: 0.129 stylish: 0.129	beauty: 0.785 smoothly: 0.516 harmonious: 0.516 passionately: 0.516 warm: 0.466 strong: 0.433 affectionately: 0.416 suited: 0.416 gentle: 0.331 confident: 0.331 qualities: 0.323 knowledgeable: 0.301 pleasant: 0.299 bright: 0.297 appealing: 0.264 lightly: 0.249 active: 0.224 enjoying: 0.216 warmth: 0.182 stylish: 0.178	secure: 1.006 comfort: 0.521 comfortable: 0.521 stylish: 0.475 cheerful: 0.458 enjoying: 0.453 suitable: 0.402 smooth: 0.351 manner: 0.331 tasting: 0.295 strong: 0.269 fluffy: 0.269 heated: 0.228 confident: 0.211 warm: 0.158 quality: 0.122 actively: 0.105 bright: 0.091 neutral: 0.076 friendly: 0.035	beauty: 0.788 manner: 0.657 pleasant: 0.525 warm: 0.382 fine: 0.310 solid: 0.273 confident: 0.258 heated: 0.233 neutral: 0.210 active: 0.191 fluffy: 0.176 thoughtful: 0.136 suitable: 0.131 bright: 0.124 enjoying: 0.098 lightly: 0.010 plain: 0.005 gentle: -0.034 stylish: -0.035 fresh: -0.045

Table 25: Top 20 warmth-related words associated with ‘woman’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
admiration: 0.929 safely: 0.797 respected: 0.621 fluffy: 0.481 solid: 0.443 committed: 0.344 friendliness: 0.286 appealing: 0.247 reliable: 0.234 quality: 0.131 enjoying: 0.127 knowledgeable: 0.106 natural: 0.043 smooth: 0.040 comfort: 0.028 rugged: 0.011 bright: -0.014 plain: -0.015 manner: -0.022 beauty: -0.035	fresh: 0.834 respect: 0.528 calmness: 0.512 enjoying: 0.390 good: 0.342 beauty: 0.260 quality: 0.249 smooth: 0.227 gentle: 0.170 solid: 0.148 secured: 0.097 warmth: 0.074 friendliness: 0.063 neutral: -0.034 plain: -0.042 bright: -0.125 attired: -0.180 kind: -0.251 strong: -0.281 confident: -0.361	skillful: 0.910 neutral: 0.620 fresh: 0.573 suitable: 0.350 respected: 0.347 safe: 0.345 safely: 0.308 appreciation: 0.266 smooth: 0.251 secure: 0.234 welcoming: 0.151 cheerful: 0.129 attitude: 0.112 lively: 0.086 delicious: 0.016 enjoying: -0.027 good: -0.047 manner: -0.108 actively: -0.131 appealing: -0.219	safe: 1.371 prayerful: 0.915 comfortable: 0.816 smooth: 0.809 kind: 0.736 appealing: 0.572 pleasing: 0.543 delicious: 0.327 manner: 0.207 neutral: 0.055 enjoying: 0.016 actively: -0.131 confident: -0.269 welcoming: -0.294 suitable: -0.307 stylish: -0.426
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
fluffy: 0.532 quality: 0.467 enthusiastic: 0.349 smooth: 0.250 neutral: 0.174 suitable: 0.161 enjoying: 0.131 manner: 0.122 cheerful: 0.096 welcoming: 0.091 plain: -0.007 bright: -0.009 natural: -0.051 actively: -0.078 fresh: -0.095 friendly: -0.131 confident: -0.339 stylish: -0.537 attitude: -0.598 warm: -0.909	calmness: 2.594 comfort: 0.816 appealing: 0.594 beauty: 0.594 enthusiastic: 0.363 rugged: 0.308 pleasant: 0.306 fine: 0.236 smooth: 0.213 suitable: 0.210 neutral: 0.180 tasting: 0.161 good: 0.158 lively: 0.149 friendly: 0.116 secure: 0.108 cheerful: 0.106 comfortable: 0.091 manner: 0.052 plain: 0.037	solid: 0.820 suitable: 0.651 smooth: 0.507 quality: 0.471 friendly: 0.457 secured: 0.321 bright: 0.312 welcoming: 0.199 natural: 0.192 kind: 0.171 plain: 0.023 stylish: 0.014 neutral: -0.051 manner: -0.121 confident: -0.176 fresh: -0.208 actively: -0.324	kind: 0.896 lively: 0.807 smooth: 0.502 fluffy: 0.498 quality: 0.483 suitable: 0.332 manner: 0.242 bright: 0.192 cheerful: 0.109 active: 0.086 good: 0.031 actively: 0.006 plain: -0.046 friendly: -0.052 natural: -0.095 welcoming: -0.111 enjoying: -0.145 fresh: -0.156 neutral: -0.203 stylish: -0.235

Table 26: Top 20 warmth-related words associated with ‘Asian’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
hot: 1.640 secure: 0.697 fine: 0.605 talented: 0.570 reliable: 0.553 intimate: 0.553 delicious: 0.546 enjoys: 0.503 appealing: 0.487 safe: 0.470 safely: 0.439 friendliness: 0.413 confident: 0.406 attired: 0.357 good: 0.355 fresh: 0.354 secured: 0.353 trust: 0.346 strong: 0.342 knowledgeable: 0.339	stylish: 0.705 manner: 0.599 warm: 0.421 confident: 0.297 enjoying: 0.289 comfortable: 0.230 respect: 0.171 passionate: 0.158 smooth: 0.113 natural: 0.055 plain: 0.039 neutral: 0.016 strong: -0.022 actively: -0.028 bright: -0.046 attired: -0.059 kind: -0.063 fresh: -0.087 warmth: -0.095 friendliness: -0.099	reliable: 1.370 prayerful: 0.827 attitude: 0.802 plain: 0.734 enthusiastic: 0.678 committed: 0.645 cheerful: 0.548 active: 0.533 stylish: 0.502 friendly: 0.497 natural: 0.463 strong: 0.463 confident: 0.427 good: 0.410 respect: 0.394 secure: 0.391 comfortable: 0.382 passionate: 0.360 knowledgeable: 0.310 quality: 0.308	pleasing: 1.157 fresh: 0.970 friendly: 0.748 taste: 0.505 enjoying: 0.450 kind: 0.282 strong: 0.233 confident: 0.196 welcoming: 0.153 neutral: 0.121 manner: 0.091 appealing: -0.383 actively: -0.433 stylish: -0.465 suitable: -0.793 delicious: -0.904
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
solid: 1.720 respect: 1.272 tasting: 1.232 cheerful: 0.597 confident: 0.432 stylish: 0.428 thoughtful: 0.393 warm: 0.204 bright: 0.200 suitable: 0.116 enthusiastic: 0.109 active: 0.009 friendly: 0.009 fresh: 0.007 plain: 0.005 actively: 0.002 neutral: 0.001 enjoying: -0.069 lively: -0.084 natural: -0.088	friendliness: 0.925 strong: 0.878 warmth: 0.866 heated: 0.755 warm: 0.675 tasting: 0.656 quality: 0.516 hot: 0.488 comfort: 0.478 stylish: 0.451 cheerful: 0.438 attitude: 0.416 safely: 0.381 enjoying: 0.351 natural: 0.339 solid: 0.323 welcoming: 0.318 lightly: 0.304 bright: 0.300 prayerful: 0.281	stylish: 0.753 thoughtful: 0.718 manner: 0.657 confident: 0.530 natural: 0.510 enjoying: 0.468 bright: 0.244 kind: 0.048 quality: 0.030 plain: 0.001 neutral: -0.011 fresh: -0.090 suitable: -0.320 smooth: -0.372 friendly: -0.579 actively: -0.798 solid: -1.009	strong: 1.033 warm: 0.852 good: 0.852 active: 0.618 bright: 0.418 confident: 0.406 stylish: 0.354 cheerful: 0.267 lively: 0.191 fresh: 0.103 neutral: 0.070 plain: 0.023 suitable: -0.099 friendly: -0.125 actively: -0.218 welcoming: -0.219 smooth: -0.352 manner: -0.424 secured: -0.508 natural: -0.885

Table 27: Top 20 warmth-related words associated with ‘Black’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
appreciation: 1.038 admiration: 0.983 good: 0.605 authenticity: 0.453 respectful: 0.246 calmness: 0.238 warmth: 0.219 respected: 0.190 bright: 0.123 respect: 0.094 kind: 0.086 gentle: 0.078 solid: 0.072 intimate: 0.066 beauty: 0.055 warm: 0.045 strong: 0.013 friendly: -0.005 confident: -0.032 comfortable: -0.062	fluffy: 1.258 comfort: 0.995 calmness: 0.810 gentle: 0.693 natural: 0.274 passionate: 0.240 solid: 0.183 warmth: 0.156 natural: 0.156 warm: 0.155 actively: 0.154 kind: 0.054 knowledgeable: 0.054 beauty: 0.033 plain: 0.031 comfortable: -0.023 respect: -0.028 attired: -0.034 friendly: -0.116 enjoying: -0.215	natural: 0.919 beauty: 0.461 plain: 0.316 safely: 0.309 knowledgeable: 0.267 safe: 0.247 comfort: 0.231 enjoys: 0.180 enthusiastic: 0.144 strong: 0.123 respect: 0.101 actively: 0.095 quality: 0.063 stylish: 0.029 enjoying: 0.021 passionate: 0.016 smooth: 0.004 good: -0.046 confident: -0.062 neutral: -0.097	prayerful: 0.908 stylish: 0.396 comfortable: 0.323 taste: 0.313 suitable: 0.100 actively: 0.013 welcoming: -0.039 neutral: -0.066 manner: -0.079 confident: -0.093 appealing: -0.275 delicious: -0.359 strong: -0.611 enjoying: -0.858
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
intimate: 1.328 beauty: 1.067 respectful: 0.825 gentle: 0.813 secured: 0.783 thoughtful: 0.652 suitable: 0.364 active: 0.295 enthusiastic: 0.105 actively: 0.081 neutral: 0.033 plain: 0.032 fresh: 0.029 warm: 0.014 welcoming: -0.083 lively: -0.087 manner: -0.112 confident: -0.231 stylish: -0.260 smooth: -0.316	respectful: 1.389 respect: 1.203 slicked: 0.851 fluffy: 0.634 tasting: 0.634 gentle: 0.453 thoughtful: 0.412 grayish: 0.352 comfortable: 0.343 kind: 0.289 pleasant: 0.207 lightly: 0.041 plain: 0.008 neutral: 0.005 beauty: -0.003 actively: -0.029 fresh: -0.050 natural: -0.086 confident: -0.109 bright: -0.136	active: 0.736 solid: 0.475 suitable: 0.318 quality: 0.290 actively: 0.221 thoughtful: 0.164 neutral: 0.135 plain: 0.018 fresh: -0.046 secured: -0.053 confident: -0.187 smooth: -0.241 manner: -0.260 kind: -0.354 neutral: -0.403 bright: -0.483 stylish: -1.134 friendly: -1.493	lightly: 0.847 kind: 0.362 quality: 0.311 active: 0.262 confident: 0.154 neutral: 0.153 warm: 0.082 smooth: 0.030 plain: 0.021 secured: -0.094 fresh: -0.123 bright: -0.124 gentle: -0.138 stylish: -0.151 natural: -0.181 actively: -0.243 manner: -0.328 fluffy: -0.451 lively: -0.520 cheerful: -1.088

Table 28: Top 20 warmth-related words associated with ‘Indian’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
slicked: 1.805 welcoming: 0.694 friendly: 0.535 trustworthiness: 0.457 reliable: 0.413 manner: 0.400 fluffy: 0.372 quality: 0.362 enjoys: 0.350 stylish: 0.342 gentle: 0.334 suitable: 0.331 delicious: 0.297 warm: 0.257 mutual: 0.234 secured: 0.225 confident: 0.219 intimate: 0.191 comfortable: 0.153 actively: 0.139	slicked: 1.277 knowledgeable: 1.037 welcoming: 0.919 good: 0.486 secured: 0.463 friendliness: 0.413 comfortable: 0.343 charity: 0.317 natural: 0.257 passionate: 0.223 warmth: 0.217 smooth: 0.178 manner: 0.126 confident: 0.094 plain: 0.094 enjoying: 0.031 neutral: -0.021 fresh: -0.022 solid: -0.031 attired: -0.051	pleasant: 1.104 warm: 0.926 talented: 0.696 comfort: 0.691 good: 0.461 lively: 0.409 respected: 0.348 appealing: 0.327 comfortable: 0.320 plain: 0.316 welcoming: 0.267 appreciation: 0.267 enjoying: 0.248 quality: 0.189 safe: 0.140 confident: 0.104 actively: 0.096 neutral: 0.078 enjoys: 0.074 delicious: 0.054	suitable: 0.696 delicious: 0.651 comfortable: 0.596 secure: 0.449 active: 0.161 strong: 0.137 enjoying: 0.109 stylish: 0.107 confident: 0.092 taste: 0.060 manner: 0.049 neutral: 0.026 smooth: -0.036 appealing: -0.120 welcoming: -0.292
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
lively: 1.046 comfortable: 0.802 lightly: 0.712 enjoying: 0.661 pleasant: 0.471 smooth: 0.445 natural: 0.364 active: 0.291 confident: 0.176 warm: 0.174 manner: 0.162 friendly: 0.128 enthusiastic: 0.101 cheerful: 0.070 stylish: 0.053 plain: 0.004 bright: -0.045 actively: -0.049 fresh: -0.053 welcoming: -0.087	thoughtfulness: 1.208 attitude: 1.036 strong: 0.967 enthusiastic: 0.655 grayish: 0.655 fluffy: 0.639 lively: 0.615 appealing: 0.586 comfortable: 0.425 safely: 0.364 enjoying: 0.334 hot: 0.301 prayerful: 0.264 active: 0.264 good: 0.215 slicked: 0.211 respect: 0.208 secured: 0.169 stylish: 0.139 warm: 0.115	enjoying: 0.730 secured: 0.500 welcoming: 0.489 actively: 0.473 friendly: 0.420 kind: 0.336 stylish: 0.282 thoughtful: 0.202 confident: 0.181 neutral: 0.053 bright: 0.029 plain: -0.004 fresh: -0.008 solid: -0.203 smooth: -0.203 manner: -0.327 quality: -0.424	pleasant: 1.106 quality: 0.541 suitable: 0.527 cheerful: 0.504 friendly: 0.434 fluffy: 0.404 gentle: 0.354 bright: 0.250 stylish: 0.229 actively: 0.193 neutral: 0.181 lively: 0.165 smooth: 0.124 manner: 0.117 confident: 0.112 good: 0.089 plain: -0.001 fresh: -0.008 warm: -0.174 secured: -0.293

Table 29: Top 20 warmth-related words associated with ‘Latino’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
attractive: 1.098 active: 1.098 trustworthiness: 0.999 prayerful: 0.818 mutual: 0.539 comfort: 0.477 authenticity: 0.468 respectful: 0.414 passionate: 0.401 respect: 0.399 actively: 0.398 gentle: 0.345 warm: 0.308 beauty: 0.298 welcoming: 0.294 suitable: 0.221 secured: 0.209 respected: 0.205 enjoying: 0.197 talented: 0.191	respectful: 1.472 fine: 1.432 bright: 0.794 quality: 0.695 actively: 0.516 fresh: 0.502 warm: 0.418 stylish: 0.373 beauty: 0.355 friendly: 0.354 solid: 0.306 kind: 0.283 passionate: 0.262 comfortable: 0.176 attired: 0.138 friendliness: 0.130 plain: 0.086 confident: 0.041 warmth: 0.008 neutral: 0.007	welcoming: 0.558 beauty: 0.399 friendly: 0.348 respect: 0.315 plain: 0.308 active: 0.259 enjoys: 0.172 smooth: 0.166 stylish: 0.145 passionate: 0.131 manner: 0.130 fresh: 0.129 actively: 0.125 lively: 0.078 appealing: 0.023 delicious: 0.009 quality: -0.082 confident: -0.110 respected: -0.120 secure: -0.135	natural: 1.214 beauty: 1.008 taste: 0.735 secure: 0.639 strong: 0.396 actively: 0.342 stylish: 0.310 welcoming: 0.191 friendly: 0.107 delicious: 0.026 manner: -0.012 confident: -0.049 neutral: -0.076 enjoying: -0.295 appealing: -0.428 suitable: -0.630
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
respectful: 1.090 gentle: 1.008 secured: 0.679 good: 0.553 enjoying: 0.515 welcoming: 0.385 warm: 0.368 smooth: 0.364 lightly: 0.356 pleasant: 0.214 fluffy: 0.191 bright: 0.176 friendly: 0.168 quality: 0.167 active: 0.160 cheerful: 0.063 actively: 0.015 plain: 0.013 fresh: 0.006 stylish: 0.001	secured: 0.932 prayerful: 0.761 fine: 0.699 secure: 0.375 solid: 0.317 rugged: 0.312 smooth: 0.276 thoughtful: 0.276 welcoming: 0.160 enthusiastic: 0.126 confident: 0.118 safely: 0.112 quality: 0.095 actively: 0.088 active: 0.083 friendly: 0.068 cheerful: 0.058 stylish: 0.041 fresh: 0.036 manner: 0.019	tasting: 1.218 good: 0.902 thoughtful: 0.178 kind: 0.174 fresh: 0.147 confident: 0.124 plain: 0.049 friendly: 0.024 neutral: -0.000 smooth: -0.075 active: -0.098 stylish: -0.120 manner: -0.181 bright: -0.247 secured: -0.261 suitable: -0.346 natural: -0.389 actively: -0.502 quality: -0.608 solid: -0.712	secured: 1.231 pleasant: 1.121 enjoying: 1.066 welcoming: 0.869 comfortable: 0.668 gentle: 0.562 natural: 0.326 manner: 0.074 quality: 0.070 plain: 0.058 fresh: 0.024 active: 0.021 friendly: -0.029 neutral: -0.033 smooth: -0.041 cheerful: -0.066 stylish: -0.097 suitable: -0.110 warm: -0.159 confident: -0.178

Table 30: Top 20 warmth-related words associated with ‘Middle Eastern’

Llama 3.2 11B	Llama 3.2 90B	InstructBLIP 7B	InstructBLIP 13B
tasting: 1.644 gratitude: 1.381 rugged: 1.016 fine: 0.861 good: 0.693 safe: 0.671 committed: 0.612 freshness: 0.612 secure: 0.575 friendliness: 0.435 natural: 0.241 quality: 0.236 attired: 0.235 calmness: 0.230 appreciation: 0.140 authenticity: 0.140 respect: 0.138 trust: 0.131 talented: 0.126 knowledgeable: 0.110	cheerful: 2.142 lively: 1.027 committed: 0.946 welcoming: 0.668 strong: 0.563 good: 0.458 natural: 0.454 friendly: 0.405 secured: 0.405 bright: 0.298 smooth: 0.290 manner: 0.279 stylish: 0.166 gentle: 0.162 attired: 0.152 knowledgeable: 0.149 beauty: 0.065 friendliness: 0.063 confident: 0.057 warm: 0.013	bright: 0.908 cheerful: 0.848 safely: 0.653 suitable: 0.405 appreciation: 0.249 committed: 0.249 secure: 0.217 talented: 0.192 enthusiastic: 0.126 manner: 0.119 comfortable: 0.109 appealing: 0.108 active: 0.097 passionate: 0.034 stylish: 0.029 natural: 0.026 delicious: 0.017 safe: 0.001 quality: 0.001 fresh: -0.062	smooth: 0.440 strong: 0.390 suitable: 0.364 appealing: 0.336 comfortable: 0.324 enjoying: 0.225 welcoming: 0.185 confident: 0.071 actively: -0.060 neutral: -0.073 stylish: -0.144 delicious: -0.220 manner: -0.310
Pixtral 12B	Pixtral Large	Qwen 2.5 7B	Qwen 2.5 72B
strong: 1.351 good: 0.824 quality: 0.469 fluffy: 0.364 manner: 0.294 natural: 0.261 comfortable: 0.205 bright: 0.179 stylish: 0.132 secured: 0.129 lightly: 0.114 fresh: 0.098 welcoming: 0.093 friendly: 0.089 actively: 0.024 confident: 0.003 plain: -0.048 neutral: -0.070 warm: -0.160 enjoying: -0.352	knowledgeable: 1.169 rugged: 0.963 hot: 0.848 solid: 0.782 heated: 0.700 fluffy: 0.600 enjoying: 0.600 safely: 0.548 suitable: 0.449 secure: 0.326 quality: 0.308 welcoming: 0.298 lightly: 0.248 manner: 0.229 prayerful: 0.226 beauty: 0.185 actively: 0.179 pleasant: 0.173 slicked: 0.173 active: 0.165	welcoming: 0.684 cheerful: 0.651 active: 0.577 natural: 0.482 actively: 0.446 friendly: 0.325 smooth: 0.184 fresh: 0.168 bright: 0.001 manner: -0.004 secured: -0.042 quality: -0.056 suitable: -0.086 plain: -0.090 neutral: -0.143 solid: -0.231 stylish: -0.446 kind: -0.566 confident: -0.782	comfortable: 1.055 strong: 0.994 natural: 0.739 fluffy: 0.543 welcoming: 0.520 friendly: 0.458 suitable: 0.377 actively: 0.299 good: 0.298 manner: 0.172 enjoying: 0.164 fresh: 0.143 gentle: 0.119 plain: -0.055 cheerful: -0.187 warm: -0.187 stylish: -0.190 neutral: -0.209 confident: -0.296 bright: -0.412

Table 31: Top 20 warmth-related words associated with ‘white (race)’