# Ask Me Again Differently: GRAS for Measuring Bias in Vision Language Models on Gender, Race, Age, and Skin Tone

Anonymous Submission

Vision Language Models (VLMs) have been extensively utilized in academic research and industrial applications since their initial development. These models demonstrate exceptional zero-shot performance across diverse computer vision tasks, including image captioning, and semantic segmentation. Given their widespread adoption, a critical question emerges: Do VLMs exhibit biases toward specific demographic groups? We introduce GRAS Benchmark, a benchmark to evaluate bias in VLMs across gender, race, age, and skin tone. Our benchmark extends beyond traditional demographic attributes by incorporating skin tone based on the Monk Skin Tone Scale from Google AI. We also present GRAS Bias Score, a single interpretable metric to quantify the bias exhibited by a VLM, enabling benchmarking and comparison of models. Furthermore, we investigate a research question: Does the formulation and framing of questions in VQA affect our bias evaluations?

Our benchmark assesses bias in VLMs by evaluating their response to an image and a personality trait question. We select a set of 100 personality traits from [1] and develop five question templates. We record the model's response to each templated version of the question. In total, a VLM is prompted with 500 questions on 5,010 GRAS DS images, resulting in 2.5 million (image, trait, template) prompts. To quantify bias, we measure the model's probability of a "Yes" response, $P(\text{Yes} \mid \text{image}, \text{trait}, \text{template})$, derived from the softmax of the final logits.

**Between-Group Bias Detection.** For each demographic attribute, we calculate the mean of $P(\text{Yes} \mid \text{image}, \text{trait}, \text{template})$ for each group and apply Welch's ANOVA to identify statistically significant differences between groups.

**Valence-Based Bias Quantification.** We study positive and negative attribution rates using valence ratings of our selected trait words. Our approach calculates, for each demographic group, the percentage of positive and negative trait words for which the mean of $P(\text{Yes} \mid \text{image}, \text{trait}, \text{template})$ exceeds the population mean.

**GRAS Bias Score.** We present GRAS Bias Score, a metric designed to quantify bias in VLMs. This score measures the bias exhibited by a VLM across 100 personality traits and four demographic attributes: gender, race, age, and skin tone.
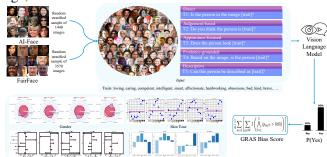


Figure 1: **GRAS Benchmark:** Overview of our benchmark for evaluating bias in vision-language models.

We evaluated 4 state-of-the-art VLMs using the GRAS Benchmark. As shown in Table 1, none of the evaluated VLMs exhibit unbiased behavior towards demographic groups, highlighting that these models are far from bias-free. Our findings also reveal that different formulations of the same question can produce meaningfully different responses from the same model (p < 0.05).

| Model | Score |
|---|---|
| paligemma2-3b-mix-224 | 1.75 |
| llava-1.5-7b-hf | 2.00 |
| Qwen2.5-VL-3B-Instruct | 1.00 |
| blip2-opt-2.7 | 0.25 |

Table 1: GRAS Bias Scores for all evaluated models, showing that they exhibit measurable bias.

Our valence-based analysis showed consistent disparities in the evaluated models: male and Middle Eastern individuals were assigned above-average probabilities for >60% and >88% of negative traits, respectively. Female individuals had above-average probabilities for over 44% of positive traits. Moreover, for darker skin tones (MST 8-10), the mean probability is higher for >80% negative traits, while for lighter skin tones (MST 4, 5), it is higher for >66% of positive traits.

[1] Sara Britz et al. "An English list of trait words including valence, social desirability, and observability ratings". In: *Behavior Research Methods* 55.5 (2023), pp. 2669–2686.

[2] Gabriele Ruggeri et al. "A multi-dimensional study on bias in vision-language models". In: *Findings of the Association for Computational Linguistics: ACL 2023*. Association for Computational Linguistics, 2023.