

Words Matter: Reducing Social Bias in AI Image Generation via Prompt Design

Junbo He, jh4737
Yuehui Ruan, yr2453

1. Synopsis

In this project, we design an experiment to assess the capability of prompt design in text-to-image generation. Compared to previous work, we introduced words other than occupation, like words in the trait and daily activity categories. We also select more contemporary and cutting-edge models into the model list. Each time, we put a word into five versions of prompts (1 baseline, 4 enhanced) and record the generation results for further analysis. We aim to provide practical insights about bias detection and mitigation so that the whole community can have a cost-effective solution to bias, a better understanding of the bias, and a wiser choice on model selection.

2. Experiment

2.1 Experiment Setup

2.1.1 Word list

We collected words across three categories: occupations, traits, and daily activities. Within each category, we grouped words based on distinct social stereotypes commonly associated with gender.

For example, some words are masculine-coded: Pilot, Firefighter, while others are feminine-coded: Nurse, secretary. Furthermore, gender neutral words are also considered: Chef, Doctor. To further assess the prompt enhancement technique, we include gender specific word pairs like Ballerina vs. Danseur, and actor vs. actress.

For each group, two or three words are selected. The full list is displayed below.

Occupation

Masculine-coded: Electrician, Firefighter, Pilot

Feminine-coded: Library, Nurse, Secretary

Gender-Neutral: Chef, Doctor

Gender-Specific: Actor, Actress, Ballerina, Danseur

Trait

Masculine-coded: Adventurous, Assertive

Feminine-coded: Sensitive, Caring

Gender-Neutral: Creative, curious

Daily Activities

Masculine-coded: Attending a business meeting, Fixing a car

Feminine-coded: Cooking dinner, Practicing yoga

Gender-Neutral: Commuting on a subway, studying at a library

2.1.2 Model list

Initially, we picked 4 models for the experiment,

Adobe Firefly Image 4, optimized for commercial use with fine-tuned safety filters and style control.

Microsoft Designer, powered by DALL-E 3, is an autoregressive diffusion model that generates highly coherent and semantically rich images.

PixArt-Alpha 1024px by Pixart produces creative and stylized images.

Stable Diffusion 2.1 by Stability AI is capable of generating realistic and high-quality images.

The first two are subscription-based models, while the other two are open source and free to use.

However, since Microsoft Designer didn't provide enough quota for the experiment, even with a subscription. We decided to eliminate it from the list.

2.1.3 Prompt Enhancement

Baseline: An image of a [occupation/trait/activity].

Inclusive I: An image of [occupation/trait/activity] who could be of any gender identity or cultural background.

Reasoning: explicit instruction of diversifying its results in identity-related generations to the model.

Inclusive II: An image of [occupation/trait/scenario] from diverse backgrounds.

Reasoning: compact and abstract instruction that gives the model freedom to generalize the idea of diversity.

Inclusive III: An image of [occupation/trait/scenario] with no assumptions about gender or ethnicity.

Reasoning: direct negation of specific stereotypes in contrast to other inclusive phrasing.

Inclusive IV: An image of [occupation/trait/activity] represented across different skin tones and genders.

Reasoning: explicit instruction of displaying variations in the generation process.

2.1.4 Experiment workflow

Data Collection

1. Select a model.
2. Select a base word.
3. Create different versions of the prompt.
4. Input to the current model and generate 4 photos using a random seed.
5. Repeat the generation if the result is extremely blurred or flawed.
6. Collect the result.

Data analysis

1. Annotate the image with gender, skin tone, and cultural markers.

2. Compute different metrics/benchmarks for words and models, respectively.
3. Retrieve some insights from both the experimental process and the computed values.
4. Aggregate and summarize the results as well as possible limitations throughout the process.

2.1.5 Annotation of the generated image

We manually assess each image based on the gender, skin tone, and cultural background of the people in it. Some of the images are blurred and ambiguous, they are also recorded for later analysis.

2.1.6 Metric

We try to obtain the following 4 metrics in our experiment.

1. Classifiability Score: It measures how many entries are discernible to provide adequate information.

$$\text{Classifiability} = \frac{\text{Number of entries with both gender and skin tone annotated}}{\text{Total number of entries}}$$

2. Diversity Score: It measures how many unique combinations of gender and skin tone appeared.

$$\text{Diversity} = \frac{\text{Number of unique (gender, skin tone) pairs}}{\text{Total number of entries}}$$

3. Synthetic IS: Mimic the true inception score, but uses softmax outputs with annotation-based heuristics. It combines the results of classifiability and diversity as a quality assessment for each image.

$$\text{Synthetic IS} = \exp \left(\frac{\log(\text{Classifiability}) + \log(\text{Diversity})}{2} \right)$$

4. CLIP Score: It measures how well the generated image matches the prompt. CLIP encodes both into vector embeddings and computes their cosine similarity.

$$\text{CLIP Score} = \frac{\vec{v}_{\text{img}} \cdot \vec{v}_{\text{text}}}{\|\vec{v}_{\text{img}}\| \cdot \|\vec{v}_{\text{text}}\|}$$

5. Identity consistency: It measures whether multiple images generated from the same prompt consistently depict the same identity.

$$\text{Identity Consistency (IC)} = \frac{1}{N} \sum_{i=1}^N \text{sim}(\vec{v}_i, \vec{v}_{\text{ref}})$$

2.2 Result

Firstly, we will display sample annotations of the generation to show how we label our result images. Each image is labeled in a similar form.

| Base Word | Prompt Version | Perceived Gender | Skin Tone | Cultural Marker |
|-------------|----------------|------------------|-----------------|---|
| Electrician | Baseline | Mixed | Light to Medium | Men and women in standard safety gear, indoor work environments |
| Electrician | Inclusive 1 | Mixed | Light to Dark | Ethnically and gender-diverse, variety in uniforms and settings |
| Electrician | Inclusive 2 | Mixed | Light to Dark | Diverse representation, cultural backgrounds visible through faces and attire |
| Electrician | Inclusive 3 | Mixed | Light to Dark | Wide range of ethnicities and gender, naturalistic and professional settings |
| Electrician | Inclusive 4 | Mixed | Light to Dark | Strong gender and racial diversity, emphasis on inclusion across roles |

We also make charts to display the distribution of images across gender, skin tone, and cultural markers.

Gender distribution

| Prompt Version | Ambiguous | Female | Male | Mixed |
|----------------|-----------|--------|------|-------|
| Baseline | 1 | 23 | 31 | 15 |
| Inclusive 1 | 0 | 39 | 1 | 30 |
| Inclusive 2 | 3 | 20 | 7 | 40 |
| Inclusive 3 | 0 | 34 | 2 | 34 |
| Inclusive 4 | 1 | 19 | 2 | 47 |

Skin tone distribution

| Prompt Version | Ambiguous | Dark | Light | Medium | Mixed |
|----------------|-----------|------|-------|--------|-------|
| Baseline | 1 | 1 | 61 | 6 | 1 |
| Inclusive 1 | 0 | 1 | 35 | 34 | 0 |
| Inclusive 2 | 3 | 2 | 27 | 34 | 4 |
| Inclusive 3 | 3 | 2 | 36 | 29 | 0 |
| Inclusive 4 | 19 | 2 | 26 | 21 | 1 |

Cultural Marker distribution

| Prompt Version | Ambiguous | Mixed | Non-Western | Western |
|----------------|-----------|-------|-------------|---------|
| Baseline | 27 | 0 | 3 | 40 |
| Inclusive 1 | 56 | 0 | 0 | 14 |
| Inclusive 2 | 50 | 2 | 1 | 17 |
| Inclusive 3 | 49 | 3 | 0 | 18 |
| Inclusive 4 | 61 | 2 | 2 | 4 |

Finally, we want to present the computed metrics as an average.

| Model | Synthetic IS (Corrected) | Classifiability | Diversity | Identity Consistency (IC) |
|----------------------|--------------------------|-----------------|-----------|---------------------------|
| Adobe FireFly | 0.813142054 | 0.87 | 0.76 | 0.455 |
| Stable Diffusion 2.1 | 0.648166645 | 0.778 | 0.54 | 0.619 |
| PixArt-Alpha | 0.471286537 | 0.667 | 0.333 | 0.556 |

| Model | Prompt Version | Avg. CLIP Score |
|----------------------|----------------|-----------------|
| Adobe Firefly | Baseline | 0.71 |
| Adobe Firefly | Inclusive 1 | 0.83 |
| Adobe Firefly | Inclusive 2 | 0.86 |
| Stable Diffusion 2.1 | Baseline | 0.68 |
| Stable Diffusion 2.1 | Inclusive 1 | 0.79 |
| Stable Diffusion 2.1 | Inclusive 2 | 0.82 |
| PixArt-Alpha | Baseline | 0.74 |
| PixArt-Alpha | Inclusive 1 | 0.84 |
| PixArt-Alpha | Inclusive 2 | 0.87 |

2.3 Analysis

Together with the annotation process and gender distribution, we found that the generation result of each word is highly skewed to one of the genders with the baseline prompt. When inclusive phrases are introduced, an interesting phenomenon happens that females are always contained or dominated in the generation results. This reminds us of an insight from our midterm experiment that inclusive phrases sometimes act like warning words to the model, and the model attempts to provide more results on socially vulnerable groups if similar phrases are present in the prompt. We still observe more mixed results, and therefore, the diversity has improved after the enhancement. Regarding skin tone distribution, we could see clear shifts from light skin tone people to an even share of light skin and medium skin tone people. For the cultural marker distribution, we experienced an extreme situation where about half resulting images are ambiguous to assess their cultural background, and this does not improve, even gets worse when inclusive prompts are employed. This could be explained by the AI model trying to introduce too many details to be realistic and diverse. Initially, they are more Western representations. After enhancement, results with additional details are displayed, making it indistinguishable regarding its cultural background. Overall, our prompt enhancement makes the distribution less skewed and more diverse.

Regarding the computed metric, for each metric, the closer the value to 1, the better the performance in generation. After computing each metric in average for the three models before and after enhancement. We found Adobe Firefly outperformed the other two models in both classifiability and diversity, hence the IS score. This is unsurprising that even with no enhancement, the model tends to provide fair results for different social groups, which is also reflected in its low IC score. This also happened in the midterm experiment, where the prompt

enhancement is less effective on Adobe Firefly, which always produces even results. Pixart, on the other hand, has the lowest diversity score, though we see remarkable improvement after the prompt design (initially very low), the overall diversity is still not satisfactory.

For the CLIP score, the closer the value is to, the better the result aligns with the prompt. We observe an increase in CLIP score for all three models since we are adding more requirements and restrictions to help AI direct their generation.

3. Discussion

3.1 Limitation

Because of the time and quota available for the experiment, a limited number of samples are generated and assessed during the process. More samples and a more comprehensive assessment may be necessary to add more authenticity to the analysis. Some metrics are computed using a small amount of data and thus can be flawed and provide little insight. The IS score in our case is mimicked through softmax outputs that could not fully explain the real capability of those models.

3.2 Research Questions

- What forms of social bias (gender, racial/skin tone, cultural) appear in outputs of text-to-image generation models when provided with neutral prompts?

In our experience during the generation process, we observe that gender bias is the most severe. In most cases, we see men dominate the results in most occupations/traits/daily activities.

- Can prompt enhancement using simple, ethically framed language significantly reduce or balance the representation of social groups in these generated images?

In our experiment, the ethically enhanced prompt does improve the representation distribution of social groups in the generated images to some extent. We see a reduction in male domination, a shift in skin tone, and an improved CLIP score.

- How consistent are the effects of such prompt modifications across different types of attribution (e.g., gender, skin tone, cultural marker)?

The gender bias is alleviated the most while skin tone only sees a clear shift from light to even share between light and medium. The cultural bias is fixed less effectively as more details are introduced, making the image indistinguishable and even vague. We conclude that prompt modification applies its effect to different extents across different types of labels.

3.3 Related works

Bias in Text-to-Image Models

Many previous works showed that biases and stereotypes in text-to-image models are reinforced frequently during the generation process.

Can be Found at “T2IAT: Measuring Valence and Stereotypical Biases in Text-to-Image Generation” by Jialu Wang, Xinyue Gabby Liu.

Prompt Engineering for Bias Mitigation

Similar experiments were done in work like “FairT2I: Mitigating Social Bias in Text-to-Image Generation via Large Language Model-Assisted Detection and Attribute Rebalancing” by Jinya Sakurai, Issei Sato. However, these are often limited to occupation categories or a narrow prompt set.

3.4 Future works

If time allowed and with paid quota, more samples should be used to assess the capability of prompt design. More variation of the prompt can also be employed if possible. With more data collected, one can create neural networks for real IS score computation. The other metrics could be more accurate and authentic in a larger dataset. A paid recognition service could also be available to help with the annotation and thus reduce the subjectivity of manual labeling.

4. Conclusion

In this paper, we introduced our experiment on examining the effectiveness of prompt enhancement in text-to-image model generation. After analyzing the result, we found the prompt modification influenced differently across the gender, skin tone, and cultural background of the generated images. Among the models we selected and assessed in the experiment, Adobe Firefly has the best performance and is affected the least by our prompt design technique. PixArt, on the other hand, needs an additional approach to improve its diversity of generation since the prompt improvement alone doesn’t guarantee the model a satisfactory performance. Overall, our technique positively affects all three models and prompt words across all categories.

5. Self-Evaluation

Junbo He

I mainly did the data collection part. I designed the collection process. My team member and I collaborate on the analysis part. I found it difficult when some of the generation results are flawed and cannot be fixed after any configuration changes. We need to decide whether to discard it or accommodate it. I also learnt how the calculated metric differed from our observation during the generation process. This could be explained by a lack of enough data or human subjectivity.

Yuehui (Ricky) Ruan

I led the analysis of image generation models in this project, focusing on representation and inclusivity across *Adobe Firefly*, *Stable Diffusion*, and *PixArt-Alpha*. I collected and coded image attributes, computed metrics like Inception Score, Inclusion Coefficient, and CLIP similarity, and conducted statistical tests to assess prompt effects. A key challenge was approximating perceptual metrics without direct model outputs, which I addressed through structured data analysis and took a huge amount of time. This project deepened my understanding of bias in generative AI and strengthened my skills in data evaluation, fairness auditing, and statistical reasoning.

6. Appendix

Github repository:

<https://github.com/JunboHe-David/ai-image-generation-6156-final-project#>

Dropbox

<https://www.dropbox.com/scl/fo/5b54brfy72f6mq0t0pik0/AFSFkt7GI3Y2Rhh4HihoTfU?rlkey=g9zskndyz578cz5l6aginq0a0&st=wpt49075&dl=0>