



# OASIS Uncovers: High-Quality T2I Models, Same Old Stereotypes

Anonymous CVPR submission

Paper ID 17

## Abstract

001 *Images generated by text-to-image (T2I) models often exhibit*  
 002 *visual biases and stereotypes of concepts such as culture and*  
 003 *profession. Existing quantitative measures of stereotypes are*  
 004 *based on statistical parity that does not align with the socio-*  
 005 *logical definition of stereotypes and, therefore, incorrectly*  
 006 *categorizes biases as stereotypes. Instead of oversimplifying*  
 007 *stereotypes as biases, we propose a quantitative measure of*  
 008 *stereotypes that aligns with its sociological definition. We*  
 009 *then propose OASIS to measure the stereotypes in a gener-*  
 010 *ated dataset and understand their origins within the T2I*  
 011 *model. OASIS includes two scores to measure stereotypes*  
 012 *from a generated image dataset: (M1) Stereotype Score*  
 013 *to measure the distributional violation of stereotypical at-*  
 014 *tributes, and (M2) WALs to measure spectral variance in*  
 015 *the images along a stereotypical attribute. OASIS also in-*  
 016 *cludes two methods to understand the origins of stereotypes*  
 017 *in T2I models: (U1) StOP to discover attributes that the*  
 018 *T2I model internally associates with a given concept, and*  
 019 *(U2) SPI to quantify the emergence of stereotypical attributes*  
 020 *in the latent space of the T2I model during image generation.*  
 021 *Despite the considerable progress in image fidelity, using*  
 022 *OASIS, we conclude that newer T2I models such as FLUX.1*  
 023 *and SDv3 contain strong stereotypical predispositions about*  
 024 *concepts and still generate images with widespread stereo-*  
 025 *typical attributes. Additionally, the quantity of stereotypes*  
 026 *worsens for nationalities with lower Internet footprints.*

## 027 1. Introduction

028 In a sociological context, stereotypes are generalized beliefs  
 029 or assumptions about a particular group of people, things,  
 030 or categories [13]. These stereotypes are widespread in the  
 031 images generated by text-to-image (T2I) models when the  
 032 input textual prompts contain concepts such as culture and  
 033 profession. For instance, consider the images in Fig. 1 gen-  
 034 erated by FLUX.1 [10], SDv3 [23], and SDv2 [43] for the  
 035 prompt “A photo of a/an <nationality> person”. There are  
 036 clear portrayals of ethnic stereotypes in attributes such as  
 037 *clothing, skin tone, and facial features* across different na-

038 tionalities, despite no references to such attributes in the  
 039 prompt. For example, the model consistently depicts an *Ira-*  
 040 *Iranian* person as a *middle-aged* or *senior* with a *long beard*,  
 041 *wearing a turban*, and dressed in *religious attire*, reinforc-  
 042 ing harmful stereotypical representations about people with  
 043 *Iranian* nationality. Besides being demographically incor-  
 044 rect, stereotypical biases in these models can lead to broader  
 045 harm. For instance, when the biased outputs of these models  
 046 are shared online, they can perpetuate damaging stereotypes  
 047 about marginalized groups, further exacerbating societal po-  
 048 larization on issues such as beauty standards, ethnicity, and  
 049 disability representation [21, 51, 55].

050 Existing methods to detect stereotypes primarily rely on  
 051 feedback from human annotators, which is both subjective  
 052 and resource-intensive. It also becomes impractical in the  
 053 era of the fast-paced development of generative models and  
 054 changing regulations. Additionally, the feedback from hu-  
 055 man annotators may be affected by their personal and politi-  
 056 cal leanings [25, 45], e.g., annotation of continuous-valued  
 057 attributes such as nose size and skin tone. Human annotation  
 058 can also affect the users’ privacy by exposing the generated  
 059 images to external evaluators.

060 In contrast, automated methods use classifiers to detect  
 061 stereotypes [16, 24, 55], overcoming several drawbacks of  
 062 human annotators. However, these methods incorrectly rely  
 063 on a general bias metric, i.e., statistical parity, as a stereotype  
 064 measure that fails to account for the directionality in the  
 065 sociological definition of stereotypes. For example, consider  
 066 a biased T2I model that generates images of predominantly  
 067 female doctors. Existing works categorize this bias as a  
 068 stereotype, although the generally known gender stereotype  
 069 associated with the concept of *doctor* is that *all doctors are*  
 070 *male* [52].

071 This paper presents a new mathematical definition of  
 072 stereotypes that aligns with the sociological definition. Build-  
 073 ing upon this formulation, we propose Open-set Assessment  
 074 of Stereotypes in Image generative models (OASIS), a novel  
 075 toolbox for quantifying stereotypes and understanding their  
 076 origins in T2I models, addressing the limitations of prior  
 077 studies. OASIS provides two metrics for measuring stereo-  
 078 types based on the distribution and spectrum of the generated

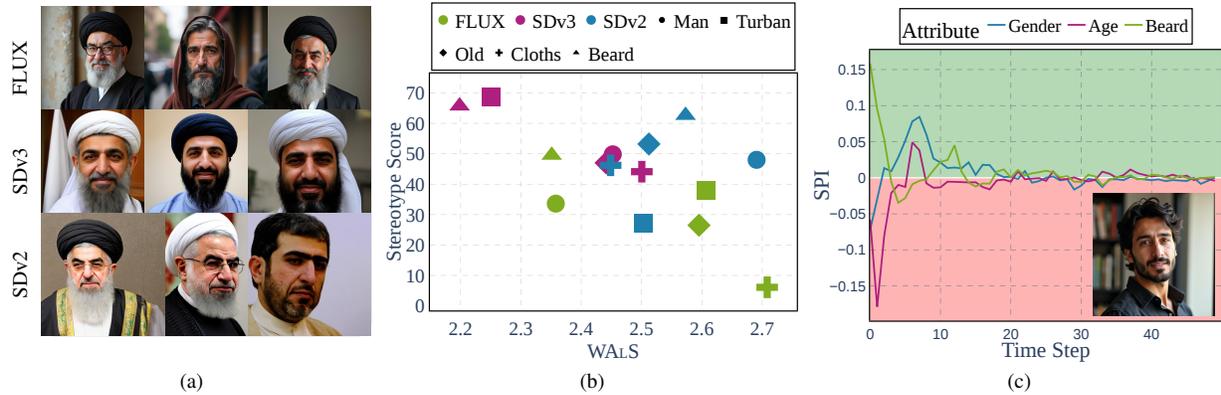


Figure 1. **Measuring Stereotypes in Text-to-Image Models.** (a) The images generated by T2I models corresponding to the prompt “A photo of an *Iranian* person” overwhelmingly contain stereotypical tropes such as *beard*, *turban*, and *religious attire* although the prompt is devoid of this information. (b) The proposed toolbox OASIS includes complementary methods for quantifying stereotypes. Stereotype Score measures the over-representation of stereotypical attributes while WAL.S measures the variance of images along these attributes. (c) SPI quantifies the emergence of stereotypes from the latent space of these models and helps understand the origin of stereotypes within a T2I model.

079 data in a feature space. OASIS comprises two additional  
 080 methods to (1) discover the stereotypical attributes that a T2I  
 081 model internally associates with a concept and (2) quantify  
 082 the emergence of stereotypical attributes in the latent space  
 083 of T2I models. Our work is an important step toward auto-  
 084 mated auditing and mitigating stereotypical content in T2I  
 085 models during development and deployment.

## 086 2. Problem Definition

087 **Definitions and Notations.** We use the term *concept* to  
 088 refer to groups of people, things, or categories related to  
 089 which stereotypes may exist, e.g., culture and profession.  
 090 We denote concepts using a random variable  $C$ . If  $C$  de-  
 091 notes the concept of nationality, then it takes values from  
 092  $\{Iranian, American, \dots\}$ . For a given concept  $C = c$ , we  
 093 define the set of potential stereotypical attributes as  $\mathcal{A}_c \subset \mathcal{A}$ ,  
 094 where  $\mathcal{A}$  is the set of all possible attributes. Every attribute  
 095  $A_i \in \mathcal{A}_c$  is a binary random variable that assumes values  
 096 from  $\{a_i^+, a_i^-\}$ , where  $a_i^+$  and  $a_i^-$  indicate the presence and  
 097 the absence of  $A_i$ , respectively. For example, if  $c = Iranian$ ,  
 098 then  $\mathcal{A}_c = \{beard, religious\ symbols, hijab, \dots\}$ . We de-  
 099 note *concepts* and *stereotypes* in different colors.

100 **Problem Setting.** The objective is to measure stereotypes  
 101 in a T2I model  $\mathcal{M}$  from the set  $\mathcal{A}_c$  that purportedly exists  
 102 related to a concept  $c$ . For example, in Fig. 1,  $c$  could corre-  
 103 spond to *Mexican nationality* and  $\mathcal{A}_c$  could include *sombrero*  
 104 and *serape*. The distribution of images  $I$  generated by  $\mathcal{M}$   
 105 conditioned on text prompt  $T(c)$  is  $p_{\mathcal{M}}(I | T(c))$ . The nota-  
 106 tion  $T(c)$  indicates that the text prompt contains information  
 107 about only the concept and not of any stereotype. To detect  
 108 the presence of  $A \in \mathcal{A}_c$ , we are provided with a dataset  $\mathcal{D}$  of  
 109  $N$  samples generated by  $\mathcal{M}$  from text prompts  $T(c)$  where  
 110  $\mathcal{D} := \{I_i | I_i \sim p_{\mathcal{M}}(I | T(c)), i = 1, \dots, N\}$ .

## 111 3. OASIS: A Stereotype Measurement and Un- 112 derstanding Toolbox

113 **Motivations.** The measurement of a stereotype related to  
 114 a concept is subjective without a formally defined metric.  
 115 Prior works have not considered the differences between  
 116 stereotypes and biases and have employed bias definitions  
 117 as stereotype metrics. The dataset  $\mathcal{D}$  is considered unbiased  
 118 w.r.t. an attribute  $A \in \mathcal{A}_c$  if

$$A | \mathcal{D} \sim \mathcal{U} \quad (1) \quad 119$$

120 where  $\mathcal{U}$  is uniform distribution. However, not all biases are  
 121 necessarily stereotypes.

122 **Quantitative Measure of Stereotype.** Stereotypes are  
 123 generalized beliefs or assumptions about a particular group  
 124 of people, things, or categories [13]. “Generalization” in  
 125 this definition can be translated to statistical terms as ex-  
 126 ceeding the true distribution of the data for a concept  $c$  in  
 127 the real world. As an example, if  $\mathcal{D}$  contains generated im-  
 128 ages of *doctors in the US* and the stereotype of interest  $A$   
 129 is *male*, the distribution of *male* in  $\mathcal{D}$  must match with its  
 130 true distribution in the real world  $P^*(A | C)$ <sup>1</sup> i.e.,  $P(A =$   
 131 *male* |  $\mathcal{D}, C = Doctor) = P^*(A = male | C = Doctor)$ .  
 132 Moreover, stereotypes are directional, which means *male*  
 133 having a smaller likelihood of *doctors in the US* compared  
 134 to the real-world distribution is not considered a stereotype,  
 135 although it is a bias. Accounting for this directionality, we  
 136 say a dataset  $\mathcal{D}$  contains stereotype  $A$  w.r.t.  $c$  if

<sup>1</sup> $P^*(A | C)$  can be obtained from census and online sources. For details, refer to § A.5.3.

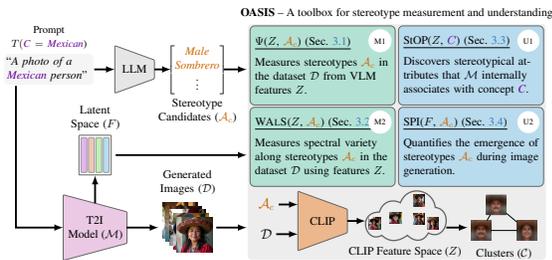


Figure 2. **An overview of OASIS.** Given a text prompt, a set of images is generated using the T2I model  $\mathcal{M}$ . Simultaneously, a stereotype candidate set is created using an LLM. OASIS then performs four quantitative analyses: (M1) Stereotype Score  $\Psi$  to measure stereotypes based on Def. 1, (M2) WALs to assess the spectral variance of  $\mathcal{D}$  w.r.t. a stereotypical attribute, (U1) StOP to discover the stereotypical attributes that  $\mathcal{M}$  internally associates with concept  $c$ , and (U2) SPI to quantify the emergence of stereotypical attributes in the latent space of  $\mathcal{M}$  during image generation.

### Definition 1. Stereotype

$$\max(0, P(A | \mathcal{D}, C) - P^*(A | C)) \geq \zeta$$

where  $\zeta$  is a margin for the violation from the real-world distribution. Note that this definition of stereotype differs from the definition of bias in Eq. (1). Our definition (i) compares the distribution of the generated dataset against its true societal distribution, and (ii) concerns the violation only along the direction of the attribute prone to be a stereotype. **Finding Stereotype Candidates.** To find open-set stereotype candidates for a concept  $c$ , we follow the approach by D’Inca et al. [21]. Let  $\mathcal{M}_{\text{LLM}}$  be a large language model (LLM). By providing prompt  $T(c)$  and a template instruction  $\mathcal{I}^2$ , we have

$$\mathcal{M}_{\text{LLM}}(T(c), \mathcal{I}) = \{(A_i, d_i^+, d_i^-) \mid i = 1, \dots, n_{\mathcal{A}_c}\} \quad (2)$$

where  $d_i^+$  and  $d_i^-$  are the descriptions for the presence and the absence of  $A_i$ , respectively. Subsequently,  $\mathcal{A}_c := \{A_1, \dots, A_{n_{\mathcal{A}_c}}\}$ . For example, let  $T(c)$  be “A photo of a doctor” and  $A_i$  be *male*.<sup>3</sup> Here,  $d_i^+$  is “A photo of a man” and  $d_i^-$  is “A photo of a woman”.

Based on these definitions, we propose OASIS, a toolbox to measure stereotypes in  $\mathcal{M}$  from distributional and spectral perspectives and to understand the origin of these stereotypical attributes in the T2I model. Given a concept  $c$ , OASIS takes in as input the dataset  $\mathcal{D}$  corresponding to a prompt  $T(c)$ , the latent space  $F$  from  $\mathcal{M}$  at every time step of image generation, and the candidate set of stereotypes  $\mathcal{A}_c$ . OASIS first extracts features  $Z$  from the images using a pre-trained vision-language model (VLM) such as CLIP [42]. Using these inputs, OASIS calculates the metrics we define below. Fig. 2 illustrates an overview of the proposed toolbox OASIS.

<sup>2</sup>Refer to § A.5.1 for more details on the template instruction.

<sup>3</sup>The number of categories for gender is restricted by the annotations of the existing datasets.

## 3.1. Stereotype Score: Measuring Stereotypes in T2I Models

Following Def. 1, stereotype score ( $\Psi$ ) of  $A \in \mathcal{A}_c$  for a given dataset  $\mathcal{D}$  and concept  $c$  is defined as

$$\Psi(A | \mathcal{D}, C) := \max(0, P(A | \mathcal{D}, C) - P^*(A | C)) \quad (3)$$

where  $P^*(A | C)$  is the real-world density of  $A$  in concept  $c$ . Using Bayes’ rule,  $P(A | \mathcal{D}, C)$  is,

$$P(A = a^+ | \mathcal{D}, C) = \frac{\prod_{i=0}^N P(A = a^+ | I_i, C)}{\sum_{a'} \prod_{i=0}^N P(A = a' | I_i, C)} \quad (4)$$

We obtain  $P(A | I_i, C)$  by means of attribute classifiers. Instead of training attribute-specific classifiers, a zero-shot predictor such as CLIP [42] can be used, where  $P(A | I_i, C)$  is obtained using a softmax over cosine similarity scores of image features and text descriptions for  $a^+$  and  $a^-$ . However, these cosine similarity scores are often numerically close [31], requiring an additional temperature parameter to obtain accurate probability measures. Therefore, in such cases, we estimate  $P(A | I_i, C)$  as

$$P(A = a^+ | I_i, C) = \mathbb{1}(\langle Z_I, Z_{a^+} \rangle_{\text{cos}} > \langle Z_I, Z_{a^-} \rangle_{\text{cos}}) \quad (5)$$

where  $\langle x, y \rangle_{\text{cos}}$  is the cosine similarity between  $x$  and  $y$ ,  $\mathbb{1}$  is the indicator function, and  $Z_I$ ,  $Z_{a^+}$ , and  $Z_{a^-}$  are features of image,  $d^+$ , and  $d^-$  from Eq. (2), respectively.

## 3.2. WALs: Measuring Spectral Variety along a Stereotype

**Motivation.** Since  $\Psi$  measures stereotypes from a distributional perspective, it is possible for a dataset  $\mathcal{D}$  to appear free of stereotypes at the cost of reduced variance along the stereotypical attribute. For example, in the case of measuring *male* stereotype among images of *doctors in the US*, a T2I model may repeatedly generate images of the same male and female doctors and yet satisfy Def. 1. Moreover, it is challenging to measure variety through human inspection due to its subjective nature, and therefore, a quantitative method to inspect variance is beneficial. To encapsulate these requirements, we propose a metric named Weighted Alignment Score (WALS) that measures the spectral alignment of the data  $\mathcal{D}$  with a given attribute  $A$ .

**Method.** To quantify the changes in a given stereotypical attribute  $A$  across images generated by a T2I model, WALs involves two steps: **1) Estimating the structure of data  $\mathcal{D}$**  through the singular value decomposition of the CLIP image features  $\mathcal{E}_I(\mathcal{D})$  i.e.,  $\mathcal{E}_I(\mathcal{D}) = U\Sigma V^T$  where  $\mathcal{E}_I$  is the image encoder of the CLIP model, **2) Finding the direction of change in  $A$** , denoted by  $\delta A$ , using one of the following two approaches: (i) estimating  $\delta A$  as the difference between the text embeddings of a pair of positive and negative descriptions,  $d^+$  and  $d^-$ ,

$$\delta A = \mathcal{E}_T(d^+) - \mathcal{E}_T(d^-), \quad (6)$$

where  $\mathcal{E}_T$  is the text encoder of the CLIP model, or (ii) estimating  $\delta A$  as the direction of maximum change along  $A$  in the image embedding space of a set of  $A$ -aware images corresponding to positive ( $a^+$ ) and negative ( $a^-$ ) categories of  $A$ , using supervised principal component analysis [5]. Detailed descriptions and proofs are mentioned in § A.6. These approaches make different assumptions, and one of these can be chosen based on the problem statement and the availability of the computational resources. The first approach assumes alignment between text and image embeddings in the CLIP model, and  $\delta A$  is more accurate when the embeddings of these modalities are more aligned. The second approach estimates  $\delta A$  accurately at the cost of increased computation due to generating two  $A$ -aware image sets and calculating the kernel matrices. Moreover, the first approach captures linear dependency, while the second one can be adopted for both linear and non-linear dependencies. We use the former approach in our experiments. Using the two components explained above, WALs measures the data variance along  $\delta A$  in the feature space, as

$$\text{WALS}(A) := \frac{\sum_{i=1}^k \sigma_i \cdot \delta A^T u_i}{\sum_{j=1}^k \sigma_j} \quad (7)$$

where  $\sigma_i$  is  $i^{\text{th}}$  singular value of  $\mathcal{D}$ , and  $u_i$  is the associated singular vector.

### 3.3. StOP: Discovering Internally Associated Stereotypical Attributes

**Motivation.** Stereotypes might occur due to T2I models internally associating a concept  $c$  with stereotypical attributes. This means that the prompts with these attributes can equivalently generate images corresponding to  $c$ . However, these attributes may not be present in  $\mathcal{A}_c$ . Therefore, qualitative methods are devised to discover these open-set attributes, which we refer to as  $\mathcal{M}$ -attributes.

**Method.** Since the distribution of stereotypical attributes is not uniform within  $\mathcal{D}$ , we have to find  $\mathcal{M}$ -attributes for individual clusters of images that share common stereotypes. Given an image dataset corresponding to concept  $c$ , we use spectral clustering [53] on CLIP features extracted from these images and visually identify clusters that share stereotypes. To discover  $\mathcal{M}$ -attributes for a given cluster with prominent stereotypes, we design a sequence optimization problem, following ZeroCLIP [49]. The solution to this optimization problem is a sequence that maximizes its mean CLIP score with the images in the chosen cluster. Formally, with a cluster  $\mathcal{D}' = \{I_1, \dots, I_n \mid 1 \leq i \leq n\}$  containing  $n$  images, the objective is

$$s^* = \arg \max_s \frac{1}{n} \sum_{i=1}^n \langle \mathcal{E}_T(s), \mathcal{E}_I(I_i) \rangle_{\cos} \quad (8)$$

where  $\mathcal{E}_I$  and  $\mathcal{E}_T$  are image and text encoders from CLIP. Following ZeroCLIP,  $s$  is produced by an LLM<sup>4</sup> that is conditioned on the starting sequence “This is a photo of”. The subsequent optimization problem reduces to iteratively finding a 2-token sequence that maximizes the mean CLIP score in Eq. (8) using beam search. Since a single prompt  $s^*$  may not contain diverse stereotypical attributes, we output the top- $K$  prompts in the final iterative step of the optimization in Eq. (8).

### 3.4. SPI: Understanding the Emergence of Stereotypes in T2I Models

**Motivation.** In addition to measuring stereotypes from generated images, it is important to quantify the aggregation of stereotypical attributes during image generation to design successful mitigation strategies. To that end, we propose stereotype propagation index (SPI) to quantify the addition of stereotypical attributes in the latent space of  $\mathcal{M}$  at each time step of image generation.

**Method.** In the flow-based models such as SDv3, the latent in each inference step is updated as  $x_{t+1} = x_t + v_{\Theta}(x_t, t, \epsilon_t)$  where  $x_t$  and  $x_{t+1}$  are the latent representation in the current and next step, respectively,  $v_{\Theta}(\cdot)$  is the velocity of  $x_t$  for time step  $t$ , and  $\epsilon_t = \epsilon_{\Theta}(x_t, t, c_p)$  is the noise predicted in time step  $t$  for latent  $x_t$  by the noise predictor  $\epsilon_{\Theta}$ , where  $c_p$  is the conditioning text prompt. The velocity decides the attributes of the generated image based on the provided text prompt. Our goal is to measure the amount of a stereotypical attribute added during each step of image generation, which requires knowing the direction of change in the attribute ( $\delta A$ ) in the latent space of the T2I model.

To find  $\delta A$  in the latent space of the T2I model, we first predict two  $A$ -aware noises that correspond to positive  $d^+$  and negative  $d^-$  descriptions of  $A$  as

$$\epsilon_t^+ = \epsilon_{\Theta}(x_t, t, d^+) \quad \epsilon_t^- = \epsilon_{\Theta}(x_t, t, d^-). \quad (9)$$

Using these predicted noises, we find the velocities that model could have in this step if the text prompt was  $A$ -aware, i.e.,  $v_{\Theta}(x_t, t, \epsilon_t^+)$  and  $v_{\Theta}(x_t, t, \epsilon_t^-)$ . Here, the direction of change in the attribute can be calculated as

$$\delta A = v_{\Theta}(x_t, t, \epsilon_t^+) - v_{\Theta}(x_t, t, \epsilon_t^-). \quad (10)$$

We define SPI as the cosine similarity between the velocity at time step  $t$  and the direction of change in the given attribute  $A$ :

$$\text{SPI}(A, t) := \langle v_{\Theta}(x_t^i, t, \epsilon_t), \delta A \rangle_{\cos} \quad (11)$$

A positive SPI means the stereotypical attribute is being added to the image in time step  $t$ , and a negative SPI means that the image is losing the stereotypical attribute  $A$ .

<sup>4</sup>We use Llama 3.1 [22]

## References

- [1] Sandhini Agarwal, Gretchen Krueger, Jack Clark, Alec Radford, Jong Wook Kim, and Miles Brundage. Evaluating clip: towards characterization of broader capabilities and downstream implications. *arXiv preprint arXiv:2108.02818*, 2021. 10
- [2] Osman Aka, Ken Burke, Alex Bauerle, Christina Greer, and Margaret Mitchell. Measuring model biases in the absence of ground truth. In *AAAI/ACM Conference on AI, Ethics, and Society*, 2021. 10
- [3] Ananya. AI image generators often give racist and sexist results: can they be fixed?, 2024. 14
- [4] Hritik Bansal, Da Yin, Masoud Monajatipoor, and Kai-Wei Chang. How well can Text-to-Image Generative Models understand Ethical Natural Language Interventions? In *Conference on Empirical Methods in Natural Language Processing*, 2022. 10
- [5] Elnaz Barshan, Ali Ghodsi, Zohreh Azimifar, and Mansoor Zolghadri Jahromi. Supervised principal component analysis: Visualization, classification and regression on subspaces and submanifolds. *Pattern Recognition*, 2011. 4, 12
- [6] Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. Easily accessible text-to-image generation amplifies demographic stereotypes at large scale. In *ACM Conference on Fairness, Accountability, and Transparency*, 2023. 10
- [7] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*, 2021. 10
- [8] Abeba Birhane, Vinay Uday Prabhu, Sanghyun Han, Vishnu Naresh Boddeti, and Sasha Luccioni. Into the LAION’s Den: Investigating Hate in Multimodal Datasets. *Advances in Neural Information Processing Systems*, 2023. 10
- [9] Abeba Birhane, Sepehr Dehdashtian, Vinay Uday Prabhu, and Vishnu Naresh Boddeti. The Dark Side of Dataset Scaling: Evaluating Racial Classification in Multimodal Models. In *ACM Conference on Fairness, Accountability, and Transparency*, 2024. 10
- [10] BlackForestLabs. FLUX. <https://blackforestlabs.ai/announcing-black-forest-labs/>, 2024. 1, 7
- [11] Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. Stereotyping Norwegian salmon: An inventory of pitfalls in fairness benchmark datasets. In *Annual Meeting of the Association for Computational Linguistics and the International Joint Conference on Natural Language Processing*, 2021. 10
- [12] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in Neural Information Processing Systems*, 2016. 10
- [13] Pedro Bordalo, Katherine Coffman, Nicola Gennaioli, and Andrei Shleifer. Stereotypes. *The Quarterly Journal of Economics*, 2016. 1, 2
- [14] Sarah E Brotherton and Sylvia I Etzel. Graduate medical education, 2022-2023. *The Journal of the American Medical Association*, 330(10):988–1011, 2023. 12
- [15] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *ACM Conference on Fairness, Accountability, and Transparency*, 2018. 10
- [16] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social biases of text-to-image generation models. In *IEEE/CVF International Conference on Computer Vision*, 2023. 1
- [17] Ching-Yao Chuang, Varun Jampani, Yuanzhen Li, Antonio Torralba, and Stefanie Jegelka. Debiasing vision-language models via biased prompts. *arXiv preprint arXiv:2302.00070*, 2023. 10
- [18] Sepehr Dehdashtian, Ruozhen He, Yi Li, Guha Balakrishnan, Nuno Vasconcelos, Vicente Ordonez, and Vishnu Naresh Boddeti. Fairness and Bias Mitigation in Computer Vision: A Survey. *arXiv preprint arXiv:2408.02464*, 2024.
- [19] Sepehr Dehdashtian, Bashir Sadeghi, and Vishnu Boddeti. Utility-Fairness Trade-Offs and How to Find Them. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [20] Sepehr Dehdashtian, Lan Wang, and Vishnu Naresh Boddeti. FairerCLIP: Debiasing CLIP’s Zero-Shot Predictions using Functions in RKHSs. *International Conference on Learning Representations*, 2024. 10
- [21] Moreno D’Incà, Elia Peruzzo, Massimiliano Mancini, DeJia Xu, Vidit Goel, Xingqian Xu, Zhangyang Wang, Humphrey Shi, and Nicu Sebe. OpenBias: Open-set Bias Detection in Text-to-Image Generative Models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 1, 3, 8, 11
- [22] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The Llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 4
- [23] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *International Conference on Machine Learning*, 2024. 1, 7

- 411 [24] Felix Friedrich, Manuel Brack, Lukas Struppek, Do- 464  
412 minik Hintersdorf, Patrick Schramowski, Sasha Luc- 465  
413 cioni, and Kristian Kersting. Fair diffusion: Instructing 466  
414 text-to-image generation models on fairness. *arXiv* 467  
415 *preprint arXiv:2302.10893*, 2023. 1, 11 468
- 416 [25] Mor Geva, Yoav Goldberg, and Jonathan Berant. Are 469  
417 We Modeling the Task or the Annotator? An Invest- 470  
418 igation of Annotator Bias in Natural Language Un- 471  
419 derstanding Datasets. In *Conference on Empirical* 472  
420 *Methods in Natural Language Processing and Interna-* 473  
421 *tional Joint Conference on Natural Language Process-* 474  
422 *ing*, 2019. 1 475
- 423 [26] Arthur Gretton, Olivier Bousquet, Alex Smola, and 476  
424 Bernhard Schölkopf. Measuring statistical dependence 477  
425 with Hilbert-Schmidt norms. In *International confer-* 478  
426 *ence on algorithmic learning theory*. Springer, 2005. 479  
427 12 480
- 428 [27] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, 481  
429 Cade Gordon, Nicholas Carlini, Rohan Taori, Achal 482  
430 Dave, Vaishaal Shankar, Hongseok Namkoong, John 483  
431 Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig 484  
432 Schmidt. OpenCLIP, 2021. 7 485
- 433 [28] Akshita Jha, Vinodkumar Prabhakaran, Remi Den- 486  
434 ton, Sarah Laszlo, Shachi Dave, Rida Qadri, Chandan 487  
435 Reddy, and Sunipa Dev. ViSAGE: A Global-Scale 488  
436 Analysis of Visual Stereotypes in Text-to-Image Gen- 489  
437 eration. In *Annual Meeting of the Association for Com-* 490  
438 *putational Linguistics*, 2024. 11 491
- 439 [29] Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, 492  
440 Samuli Laine, Timo Aila, and Jaakko Lehtinen. Apply- 493  
441 ing guidance in a limited interval improves sample and 494  
442 distribution quality in diffusion models. *arXiv preprint* 495  
443 *arXiv:2404.07724*, 2024. 10 496
- 444 [30] Marie Lamensch. Generative AI tools are perpetuating 497  
445 harmful gender stereotypes. *Centre for International* 498  
446 *Governance Innovation*, 14, 2023. 14 499
- 447 [31] Victor Weixin Liang, Yuhui Zhang, Yongchan Kwon, 500  
448 Serena Yeung, and James Y Zou. Mind the gap: Under- 501  
449 standing the modality gap in multi-modal contrastive 502  
450 representation learning. In *Advances in Neural Infor-* 503  
451 *mation Processing Systems*, 2022. 3 504
- 452 [32] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou 505  
453 Tang. Deep learning face attributes in the wild. In 506  
454 *IEEE/CVF International Conference on Computer Vi-* 507  
455 *sion*, 2015. 14 508
- 456 [33] Alexandra Sasha Luccioni and Joseph D Viviano. 509  
457 What’s in the box? a preliminary analysis of unde- 510  
458 sirable content in the common crawl corpus. *arXiv* 511  
459 *preprint arXiv:2105.02732*, 2021. 10 512
- 460 [34] Sasha Luccioni, Christopher Akiki, Margaret Mitchell, 513  
461 and Yacine Jernite. Stable bias: Evaluating societal rep- 514  
462 resentations in diffusion models. *Advances in Neural* 515  
463 *Information Processing Systems*, 2024. 11 516
- [35] Helmut Lütkepohl. *Handbook of matrices*. John Wiley & Sons, 1997. 12 464
- [36] Ranjita Naik and Besmira Nushi. Social biases through the text-to-image generation lens. In *AAAI/ACM Conference on AI, Ethics, and Society*, 2023. 10 466
- [37] Leonardo Nicoletti and Dina Bass. "Humans are Biased. Generative AI is Even Worse", 2023. 14 468
- [38] OpenAI. Chatgpt 4o. <https://chat.openai.com>, 2024. 7 469
- [39] OpenAI. Chatgpt o1-preview. <https://chat.openai.com>, 2024. 7 470
- [40] Rishubh Parihar, Abhijnya Bhat, Abhipsa Basu, Saswat Mallick, Jogendra Nath Kundu, and R Venkatesh Babu. Balancing Act: Distribution-Guided Debiasing in Diffusion Models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 11 471
- [41] Hoang Phan, Andrew Gordon Wilson, and Qi Lei. Controllable Prompt Tuning For Balancing Group Distributional Robustness. *arXiv preprint arXiv:2403.02695*, 2024. 10, 11 472
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning Transferable Visual Models from Natural Language Supervision. In *International Conference on Machine Learning*, 2021. 3, 14 473
- [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE/CVF conference on computer vision and pattern recognition*, 2022. 1, 7 474
- [44] Bashir Sadeghi, Sepehr Dehdashtian, and Vishnu Bodeti. On Characterizing the Trade-off in Invariant Representation Learning. *Transactions on Machine Learning Research*, 2022. Featured Certification. 10 475
- [45] Maarten Sap, Swabha Swayamdipta, Laura Vianna, Xuhui Zhou, Yejin Choi, and Noah A Smith. Annotators with Attitudes: How Annotator Beliefs And Identities Bias Toxic Language Detection. In *North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022. 1 476
- [46] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 10 477
- [47] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. In *Advances in Neural Information Processing Systems*, 2022. 7 478

- 517 [48] Preethi Seshadri, Sameer Singh, and Yanai Elazar. The  
518 Bias Amplification Paradox in Text-to-Image Genera-  
519 tion. In *Conference of the North American Chapter of*  
520 *the Association for Computational Linguistics: Human*  
521 *Language Technologies*, 2024. 10
- 522 [49] Yoav Tewel, Yoav Shalev, Idan Schwartz, and Lior  
523 Wolf. Zerocap: Zero-shot image-to-text generation for  
524 visual-semantic arithmetic. In *IEEE/CVF Conference*  
525 *on Computer Vision and Pattern Recognition*, 2022. 4
- 526 [50] Eddie Ungless, Björn Ross, and Anne Lauscher. Stereo-  
527 types and Smut: The (Mis) representation of Non-  
528 cisgender Identities by Text-to-Image Models. In *Find-*  
529 *ings of the Association for Computational Linguistics*,  
530 2023. 10
- 531 [51] Adriana Fernández de Caleyá Vázquez and Eduardo C  
532 Garrido-Merchán. A Taxonomy of the Biases of the  
533 Images created by Generative Artificial Intelligence.  
534 *arXiv preprint arXiv:2407.01556*, 2024. 1
- 535 [52] Lauren Vogel. When people hear "doctor", most still  
536 picture a man. *CMAJ: Canadian Medical Association*  
537 *journal= journal de l'Association medicale canadi-*  
538 *enne*, 191(10):E295–E296, 2019. 1
- 539 [53] Ulrike Von Luxburg. A tutorial on spectral clustering.  
540 *Statistics and computing*, 2007. 4
- 541 [54] WorldPopulationR. "Internet Users by Country", 2024.  
542 [Online; accessed 1-Oct-2024]. 8
- 543 [55] Cheng Zhang, Xuanbai Chen, Siqi Chai, Chen Henry  
544 Wu, Dmitry Lagun, Thabo Beeler, and Fernando De la  
545 Torre. Iti-gen: Inclusive text-to-image generation. In  
546 *IEEE/CVF International Conference on Computer Vi-*  
547 *sion*, 2023. 1, 10, 11

## A. Appendix

In our main paper, we proposed OASIS for quantifying stereotypes and understanding their origins in T2I models. Here, we provide some additional analysis to support our main results. The appendix section is structured as follows:

1. Results in § A.1
2. Related Work in § A.3
3. Concluding Remarks in § A.4
4. Implementation Details in § A.5
5. Finding  $\delta A$  Using  $A$ -Aware Generated Images in § A.6
6. More Results on SPI and average SPI in § A.7
7. More Results on T2I models' Stereotypical Predispositions in § A.8
8. Limitations in § A.9
9. Importance of Detecting Stereotypes in § A.10
10. Qualitative Descriptions of the Generated Datasets in § A.11

### A.1. What does OASIS Uncover about Stereotypes in T2I Models?

We apply OASIS on three open-weight T2I models – SDv2 [43], SDv3 [23], and FLUX.1<sub>[dev]</sub> [10]. In the first step, as illustrated in Fig. 2, we generate a dataset of 2000 images of people from each of the *nationalities* and with each T2I model. In the next step, for each *nationality*, a candidate set for stereotypes and their descriptions is generated according to Eq. (2). We used ChatGPT o1-preview [39] and ChatGPT 4o [38] as  $\mathcal{M}_{LLM}$  in Eq. (2). Implementation details are mentioned in § A.5.

#### A.1.1. Lower, Yet Significant Stereotypes in Newer T2I Models

Table 1. **Stereotype Score.** Comparison of three T2I models, SDv2, SDv3, and FLUX.1 on stereotype score in three *nationalities*.  $P^*(A | C)$  is the true density of the attribute obtained from real-world statistics (details provided in § A.5.3),  $P(A | \mathcal{D}, C)$  is the density of the attribute in the generated dataset, and  $\Psi(A | \mathcal{D}, C)$  is the stereotype score. All values are in %.

Stereotype Candidate	$P^*(A   C)$	SDv2		SDv3		FLUX.1 <sub>[dev]</sub>		
		$P(A   \mathcal{D}, C)$	$\Psi(A   \mathcal{D}, C)$	$P(A   \mathcal{D}, C)$	$\Psi(A   \mathcal{D}, C)$	$P(A   \mathcal{D}, C)$	$\Psi(A   \mathcal{D}, C)$	
<i>Iranian</i>	Man	50	98	48	99.8	49.8	83.6	33.6
	Wearing Turban	0.2	27.3	27.1	69	68.8	38.2	38
	Old	40	93.2	53.2	87	47	66.5	26.5
	Traditional Cloths	50	96.2	46.2	94.1	44.1	56.1	61
	Beard	34	96.6	62.6	99.7	65.7	83.5	49.5
<i>Indian</i>	Man	51	78.5	27.5	78.1	27.1	31.6	0
	Turban	2	2.2	0.2	0.9	0	0.1	0
	Mustache	25	17.7	0	12.4	0	25.9	0.9
	Tilak/Bindi	50	61.7	11.7	59.3	9.3	86.7	36.7
	Vibrant Color Cloths	50	41.5	0	58.3	8.3	53.8	3.8
<i>Mexican</i>	Man	48	95.1	47.1	85	37	50.1	2.1
	Hat	50	77.3	22.3	49.2	0	94.4	44.4
	Sombrero	50	56.6	6.6	17.6	0	58.6	8.6
	Mustache	25	77.8	52.8	34.1	9.1	84.7	59.7
	Embroidered Clothing	50	82.6	32.6	45.9	0	94.2	44.2

We use CLIP ViT-G-14 from OpenCLIP [27] trained on LAION2B [47] to estimate  $P(A | \mathcal{D}, C)$ . Table 1 compares the T2I models in terms of their stereotype scores defined in Sec. 3.1 from the images generated by these models corresponding to three *nationalities* – *Iranian*, *Indian*, and *Mexican*. Although the fidelity of the generated images has improved dramatically from SDv2 to SDv3 and

585 FLUX.1, our results demonstrate that stereotype scores of  
 586 newer models are generally lower than those of the older  
 587 ones. However, in some cases, there are exceptions. As  
 588 an example, when generating images of *Mexican person*,  
 589 FLUX.1 depicts 84.7% of the faces with *mustache* while  
 590 SDv2 and SDv3 generate 77.8% and 34.1% faces with *mus-*  
 591 *tache*, respectively. In high-level attributes such as gender,  
 592 FLUX.1 has lower stereotype scores than other models. For  
 593 example, in the case of *Iranian*, FLUX.1 depicts 83.6% of  
 594 images as *man*. But in comparison, 98% and 99.8% of the  
 595 images generated by SDv2 and SDv3, respectively, depict  
 596 *man*.

**Remark.** Existing bias definitions are not applicable for some attributes studied in Tab. 1. E.g., a T2I model needs to depict 50% of the images of *Iranian* with *turban* to be unbiased according to Eq. (1), which incorrectly represents *Iranian people* among whom only 0.2% wear *turban*.

597

Table 2.  $P(A = \textit{man} \mid C, D)$  for  $C = \textit{doctor}$ ,  $C = \textit{Iranian doctor}$ , and  $\textit{Indian Doctor}$ .

Model	<i>Doctor</i>	<i>Indian Doctor</i>	<i>Iranian Doctor</i>
SDv2	93	97 (+4)	98 (+5)
SDv3	78	98 (+20)	100 (+22)
FLUX.1	93	100 (+7)	100 (+7)

598 Previous works have noted the gender imbalance in the  
 599 generated images for certain professions such as *doctors* and  
 600 *teachers* [21]. We observe a similar trend in the newer T2I  
 601 models as shown in Table 2. However, SDv3 has a lower  
 602 gender imbalance compared to SDv2 for *doctor*. We hypoth-  
 603 esize that this is due to the data balancing methods taken  
 604 to ensure unbiased gender representation in the images of  
 605 *doctor* following the scrutiny it has faced. However, the  
 606 imbalance worsens when a *nationality* is added to the profes-  
 607 sion (e.g., *Iranian doctor*). This example demonstrates that  
 608 stereotype mitigation through data balancing is insufficient  
 609 against intersectional stereotypes as it is infeasible to collect  
 610 data samples corresponding to every possible combination.

## 611 A.2. Images of Under-Represented Nationalities 612 Contain More Stereotypes

613 T2I models are often trained on image-caption pairs that are  
 614 scraped from the Internet. Therefore, their training data may  
 615 be biased by the Internet footprint of various nationalities.  
 616 To investigate the impact of a nationality’s Internet footprint  
 617 on stereotypes in T2I models, we compare the stereotype  
 618 scores of generated images from various nationalities against  
 619 their corresponding number of Internet users. We consider  
 620 generated images corresponding to *Indian*, *Mexican*, and  
 621 *Iranian* nationalities, which have populations of 881.3 mil-

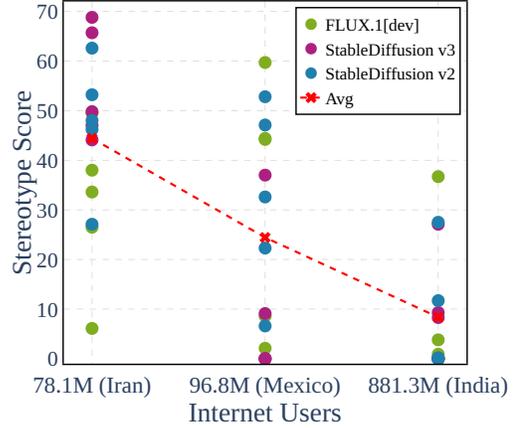


Figure 3. Comparing stereotype scores for nationalities against the number of Internet users shows that stereotypes are higher for under-represented nationalities.

622 lion, 96.8 million, and 78.1 million, respectively [54]. Fig. 3  
 623 presents the stereotype scores across different attributes for  
 624 each country and model. The results indicate that the max-  
 625 imum and the average stereotype scores for a nationality  
 626 decrease as the number of Internet users increases. These  
 627 findings suggest that the stereotypes in T2I models may be  
 628 exacerbated for under-represented nationalities when trained  
 629 on image-caption pairs from the Internet.

### 630 A.2.1. Effective T2I Model Comparison Requires Both 631 Stereotype Score and WALs

632 Fig. 4 compares the WALs for T2I models on three nation-  
 633 alities. A higher  $WALS(A)$  indicates more variance in the  
 634 images along the attribute  $A$ . We observe that FLUX.1 gen-  
 635 erates images that show a higher variety in clothing items  
 636 such as *hats* and *turbans*, but have a lower variance regard-  
 637 ing facial attributes such as *beard* and *mustache* across all three  
 638 nationalities. In contrast, images generated by SDv2 show  
 639 a higher variance on *beard* and *mustache* than on *clothing-*  
 640 *related* attributes.

641 As mentioned in Sec. 3.2, stereotype score and WALs are  
 642 complementary measures of stereotypes. We can compare  
 643 the models jointly on these scores to verify if models demon-  
 644 strate lower stereotypes at the cost of lower variety. Fig. 5  
 645 plots stereotype score and WALs for various T2I models  
 646 and attributes for each *nationality*. An ideal T2I model must  
 647 have a low stereotype score and a high WALs and there-  
 648 fore must appear towards the bottom-right corner of these  
 649 plots. We observe that some models have lower stereotypes  
 650 while having lower attribute variance. For instance, images  
 651 from SDv3 tend to have lower WALs across all *nationali-*  
 652 *ties*, although they succeed in reducing stereotypes in some  
 653 attributes. These observations highlight the importance of  
 654 employing both distributional (stereotype score) and spec-  
 655 tral (WALS) metrics together to compare the T2I models.

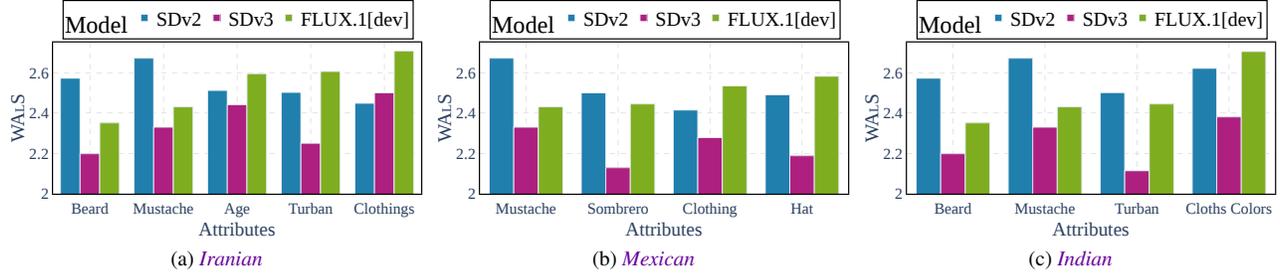


Figure 4. WALs: Comparison of SDv2, SDv3, and FLUX.1 on spectral variance in the generated images across different attributes, calculated for *Iranian*, *Mexican*, and *Indian* nationalities.

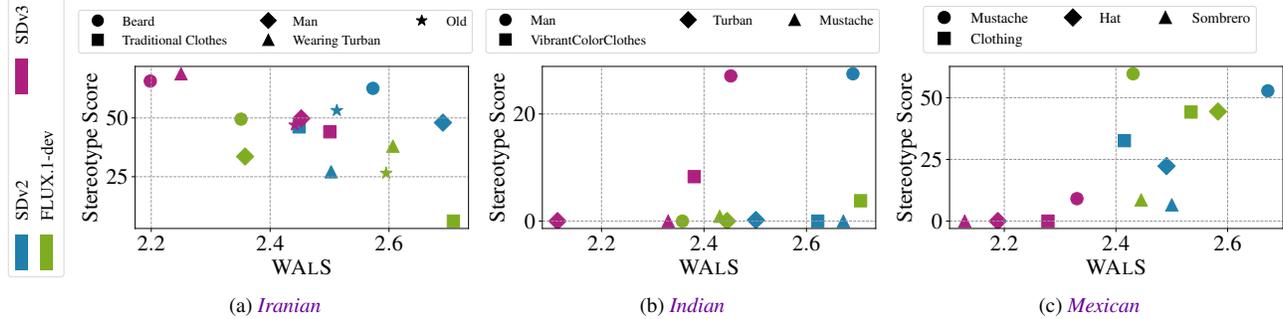


Figure 5. Comparison of T2I models based on stereotype scores and WALs for three nationalities. Different colors show different T2I models and the shapes of the markers denote the *attributes*.

656  
657

### A.2.2. T2I Models Internally Associate Concepts with Stereotypes

Table 3. **StOP** first identifies image clusters for each concept using spectral clustering. The averages of the images from these clusters are shown in the second column. **StOP** finds the captions shown in the third column by solving the optimization problem in Eq. (8). These captions contain stereotypical attributes such as “*Imam*” and “*brero*”. The fourth column shows the images generated using these optimized captions. Unsurprisingly, these images contain insignia of the corresponding culture.

Culture	Cluster average	Optimized prompts	Samples from highlighted prompt
<i>Iranian</i>		“This is a photo of \u093f\u092e Imam” “This is a photo of \u0935 reb” “This is a photo of \u093f\u0935 Sheikh”	
<i>Mexican</i>		“This is a photo ofbrero mayoc” “This is a photo ofbrero Garcia” “This is a photo ofbrero pastor”	
<i>American</i>		“This is a photo of EO Democrat” “This is a photo of: border counselor” “This is a photo of: border ambassador”	

658  
659  
660  
661  
662  
663  
664  
665  
666

We use **StOP** to discover the internal associations that the T2I model  $\mathcal{M}$  makes with a given concept  $c$ . In Table 3, we show  $\mathcal{M}$ -attributes in FLUX.1 discovered using **StOP** for three concepts: *Iranian*, *Mexican*, and *American*. We obtain clusters of images using spectral clustering on the CLIP features of aligned faces and manually identify those with shared stereotypes. The average of the faces in the clusters are shown in the second column. The attributes that we expect **StOP** to discover can be visually identified

from these averaged images. For example, the average of the cluster corresponding to *Iranian* shows an *old man* wearing a *turban* and sporting a *long beard*, characteristic of the Islamic religious leaders in Iran. Therefore, the expected  $\mathcal{M}$ -attributes include religious terminology. In the third column, we show some of the optimized prompts that **StOP** produces. The optimized prompts contain Unicode characters in vernacular languages. The optimized prompts corresponding to *Iranian* images include religious terms such as “*Imam*” and “*Sheikh*”. Similarly, the optimized prompts for *Mexican* images include “*brero*” (short for *sombrero*). In the last column, we input one of these prompts to FLUX.1 to visually inspect the resulting images. Unsurprisingly, the images generated from these optimized prompts are visually similar to those generated from prompts containing only *nationality*. For example, *the US national flag* can be seen as a blurred background in the cluster average and is also present in the samples generated from optimized prompts for *American person*.

667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685

### A.2.3. Stereotypical Attributes Emerge in the Early Steps of Image Generation

We quantify the emergence of stereotypical attributes during image generation in FLUX.1 and SDv3 for image prompts of the form “A photo of an  $\langle$ nationality $\rangle$  person” using SPI. For a given stereotype, we first obtain positive and negative descriptions corresponding to it. For example, for the attribute *age*, “old” and “young” were used in  $d^+$  and  $d^-$

686  
687  
688  
689  
690  
691  
692  
693

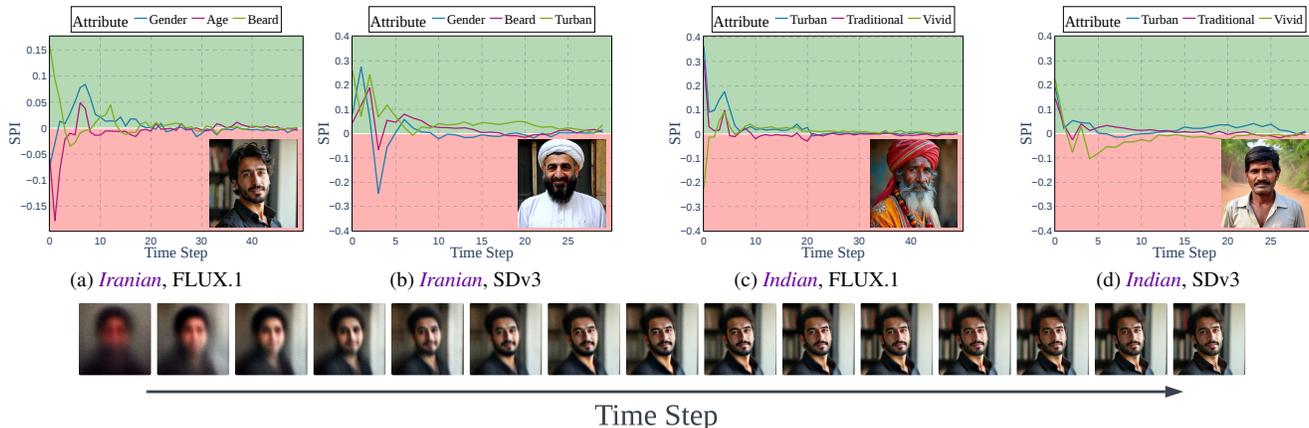


Figure 6. **SPI** tracks the change in attributes in the image generation processes of FLUX.1 and SDv3. We observe that these attributes are affected during the early time steps of image generation.

694 in Eq. (2), respectively. SPI is then calculated as the cosine  
 695 similarity between  $\delta A$  and velocity of the latent  $x_t$  at time  
 696 step  $t$  as shown in Eq. (11). We plot  $\text{SPI}(A, t)$  for all time  
 697 steps during the generation of four images from *Iranian* and  
 698 *Indian* nationalities in Fig. 6. We observe that a relatively  
 699 high amount of information on attributes such as *beard*,  
 700 *traditional cloths*, and *sombrero* is added to the images at  
 701 the first step of generation in FLUX.1 indicating that the  
 702 model readily associates these attributes with the concept.  
 703 Specifically, we observe that the stereotypical attributes arise  
 704 during the earlier time steps of generation in both *Iranian* and  
 705 *Indian* images. In the example of *Iranian* person in Fig. 6a,  
 706 we observe that *age* and *beard* attributes form in the image  
 707 within the first 3 time steps and *gender* attribute emerges at  
 708 time step 7. After time step 20, the changes in these attributes  
 709 are negligible. Similarly, stereotypical attributes form within  
 710 the first 20 time steps of generating an image from *Indian*  
 711 nationality and undergo little change afterward. Additional  
 712 results for other nationalities are provided in § A.7.

#### 713 A.2.4. T2I Models have Stereotypical Predispositions 714 about Concepts

715 In § A.2.3, we noted that stereotypical attributes aggregate in  
 716 the early steps of image generation. A question that naturally  
 717 follows this observation is: *are T2I models predisposed to*  
 718 *generate stereotypical images for a given concept?* This can  
 719 be answered by considering the velocity  $v_{\Theta}(x_t, t, \epsilon_t)$  of the  
 720 early time steps since they guide towards the mean of the  
 721 data

722 distribution [29]. This enables us to identify the stereo-  
 723 typical predispositions qualitatively. For each time step  $t$ ,  
 724 we estimate the final time step image  $\hat{x}_T$  based on velocity  
 725  $v_{\Theta}(x_t, t, \epsilon_t)$  as  $\hat{x}_T = x_t + v_{\Theta}(x_t, t, \epsilon_t)(T - t)$ , as illustrated  
 726 in Fig. 7a. Fig. 7b shows these images for three samples  
 727 corresponding to *Iranian person*. The images generated us-

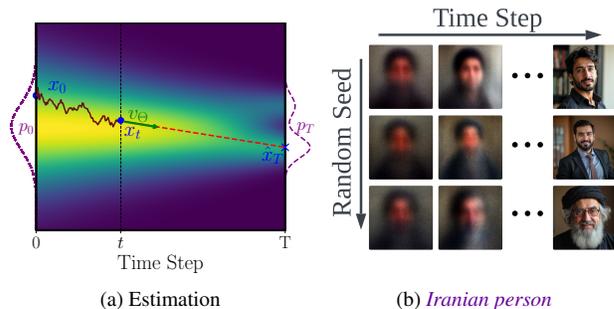


Figure 7. Stereotypical predisposition in T2I models for *Iranian person*.

ing the velocity at time  $t = 0$  appear to be of a person with  
 728 *turban* and *beard*, even when the final generated images  
 729 lack these attributes. Conflating with our observations from  
 730 SS A.2.2 and A.2.3, we conclude that T2I models associate  
 731 stereotypical attributes with seemingly innocuous prompts.  
 732 Additional results are provided in § A.8.  
 733

#### 734 A.3. Related Work

735 Many studies have shown that deep learning models tend  
 736 to learn and, at times, amplify the biases present in their  
 737 datasets [1, 2, 8, 9, 11, 12, 15, 17–20, 33, 41, 44], and T2I  
 738 models are no exception. Most of the existing work about  
 739 stereotypes in T2I models has focused on gender and ethnic  
 740 biases in the generated images. Some studies have shown  
 741 that prompts play a significant role in the bias generated by  
 742 T2I models [4, 48, 55]. Seemingly neutral prompts lead to  
 743 geographical biases favoring Western nations such as the  
 744 US and Germany, leading to lighter skin tones and Western  
 745 norms in the images [6, 36], while prompts containing cer-  
 746 tain cultural and gender terms sometimes generate NSFW im-  
 747 ages, reflecting the biases in the training datasets [7, 46, 50].

748 Luccioni et al. [34] measured distributional biases in profes-  
749 sions w.r.t a closed set of genders and ethnicities.

750 Unlike these works, we use an open set of stereotypes ob-  
751 tained from an LLM, following [21]. Although these studies  
752 have achieved breadth in terms of sources for stereotypes,  
753 they have primarily used statistical parity as the definition of  
754 stereotype. For example, [28] uses “stereotype tendency” de-  
755 fined as the ratio of the likelihood of a stereotype appearing  
756 in a group to that of it appearing in the general population,  
757 ignoring the directionality of stereotypes. In contrast, we  
758 measure stereotypes following their true sociological defini-  
759 tion. We additionally provide insights into the origins of the  
760 stereotypical attributes in T2I models.

#### 761 A.4. Concluding Remarks

762 This paper proposed OASIS to measure and understand the  
763 origin of stereotypes in T2I models based on a quantitative  
764 measure that aligns with the sociological definition of stereo-  
765 type. OASIS includes: (M1) Stereotype Score (Sec. 3.1) to  
766 measure the directional violation of the true stereotypical at-  
767 tribute distribution in the T2I model, (M2) WALs (Sec. 3.2)  
768 to measure the spectral variety of the generated images along  
769 the stereotypical attributes, (U1) StOP (Sec. 3.3) to discover  
770 the stereotypical attributes that the T2I model internally asso-  
771 ciates with a concept, and (U2) SPI (Sec. 3.4) to measure the  
772 emergence of stereotypical attributes during image genera-  
773 tion from the latent space. Despite the considerable progress  
774 in the image fidelity of T2I models, using OASIS, we con-  
775 clude that the newer models such as FLUX.1 and SDv3 have  
776 strong stereotypical predispositions about concepts and still  
777 struggle to avoid stereotypical attributes in the images.

778 **Recommendations.** OASIS unveils the extent of stereo-  
779 types in T2I models. However, commonly pursued solutions  
780 for correcting biases in generative models such as data bal-  
781 ancing are not suitable for resolving stereotypes due to the  
782 sheer number of concepts that could potentially have stereo-  
783 types. Additionally, concepts such as *nationalities* worsen  
784 stereotypes in unrelated concepts such as *doctors* as ob-  
785 served in Tab. 2. It is infeasible to collect data samples at  
786 the intersection of multiple concepts. Therefore, training-  
787 time mitigation and post hoc correction techniques that are  
788 tailored to remove stereotypes in T2I models must be de-  
789 veloped[24, 40, 41, 55]. Our observations also underscore  
790 the need for increased participation of under-represented  
791 communities in the development of large generative models.

#### 792 A.5. Implementation Details

793 In this section we provide all the details that are employed  
794 in the experiment section in order to obtain the mentioned  
795 results.

##### 796 A.5.1. LLM Template Instruction

797 To generate an open-set candidate for stereotypes correspond-  
798 ing to a given *concept*, we follow the approach provided in

[21]. The template instruction  $\mathcal{I}$  that was used in Eq. (2) is  
as follows:

I am using a text-to-image model to generate facial images of individuals from various nationalities. These images depict faces along with a surrounding margin that may include headwear or clothing visible above the chest. I need a candidate set of at least 15 potential stereotypes to monitor for in the generated images. I will provide the prompts used to generate the images, and I would like you to identify this candidate set. Additionally, I plan to check the presence or absence of these attributes in the images using a CLIP model. For each attribute, I need at least two sentences: one that positively describes the existence of the attribute and another that describes its absence.

An example of the output produced by ChatGPT o1-  
preview is,

Certainly! Based on your prompt, “A *photo of a Tunisian person*,” here is a list of potential attributes to look for in the generated images. These attributes are associated with cultural elements and common depictions of Tunisian people. For each attribute, I’ve provided a positive and a negative sentence that can be used with a CLIP model for zero-shot classification.

1. Wearing a Traditional Chechia (Hat)
  - Positive: "The person is wearing a traditional Tunisian chechia."
  - Negative: "The person is not wearing any headwear."
2. Traditional Clothing
  - Positive: "The person is wearing traditional Tunisian clothing."
  - Negative: "The person is wearing modern Western clothing."
3. Wearing a Djebba (for men)
  - Positive: "The man is wearing a traditional djebba."
  - Negative: "The man is not wearing a djebba."
- ⋮
15. Wearing a Sefsari (for women)
  - Positive: "The woman is wearing a traditional sefsari."
  - Negative: "The woman is not wearing a sefsari."

### 805 A.5.2. Image Generation Prompts

806 To generate images from FLUX.1 and SDv3, we used the  
807 prompt template “A photo of <nationality> person”, where  
808 <nationality> could take values like *American*, *Indian*, etc.  
809 For SDv2, we observed that the images could be of lower  
810 quality or contain the torso of the person. Therefore, for  
811 each images, we sampled a prompt randomly from the fol-  
812 lowing array of prompts: “A photo of <nationality> person”,  
813 “A picture of <nationality> person”, “A portrait photo of  
814 <nationality> person”, “A front profile photo of <nationality>  
815 person”.

### 816 A.5.3. Obtaining True Distributions

817 Stereotype score is measured as the violation of the true un-  
818 derlying distribution of an attribute given a concept, denoted  
819 by  $P^*(A | C)$ , in the generated images. One could obtain  
820  $P^*(A | C)$  from official census data and online statistics. For  
821 example, Brotherton and Etzel [14] provides various demo-  
822 graphic details about doctors in the US such as ethnicity and  
823 gender in various specializations. For attributes where it is  
824 difficult to obtain precise statistics, e.g., traditional cloth-  
825 ing, we consider their presence a choice and assign a 50%  
826 chance for their presence. For example, for *mustache* for peo-  
827 ple from *Mexican* nationality, we calculate  $P^*(\text{mustache} |$   
828 *Mexican}) = 0.5 \times P^\*(\text{male} | \text{Mexican}) \approx 0.255.*

## 829 A.6. Finding $\delta A$ Using $A$ -Aware Generated Images

830 As mentioned in Sec. 3.2, to find the direction of change in  
831  $A$ , we propose two approaches: (i) using text embeddings of  
832 a pair of positive and negative descriptions,  $d^+$  and  $d^-$ , and  
833 (ii) using  $A$ -aware generated images. The first approach is  
834 explained in Sec. 3.2 and in this section, we explain how to  
835 find  $\delta A$  using  $A$ -aware generated images in both linear and  
836 non-linear cases.

### 837 A.6.1. Linear $\delta A$

838 As mentioned in Sec. 3.2, using  $A$ -aware generated images  
839 to find  $\delta A$  can be more precise than using text embeddings.  
840 As an example, for finding the direction of change in *male*  
841 for images corresponding to “A photo of an *Iranian* person”,  
842 two sets of images using prompts “A photo of a *man*” and  
843 “A photo of a *woman*” are created. The set of CLIP features  
844 of these images are denoted by  $Z_A = \{z_i\}_{i=1}^m$  and their  
845 corresponding labels of  $A$  as  $Y_A = \{y_i\}_{i=1}^m$ . Then we find  
846 an orthogonal transformation matrix  $\Gamma$  that maps  $Z_A$  to a  
847 subspace that maximizes the variance of the labeled data  
848 using supervised principal component analysis [5] that max-  
849 imizes the dependency between the mapped data  $\Gamma^T Z_A$  and  
850  $Y_A$ . Hilbert-Schmidt Independence Criterion (HSIC) [26]  
851 is employed as the dependence metric where its empirical  
852 version is defined as  $\text{HSIC}^{\text{emp}} = \text{Tr}\{H K_{ZZ} H K_Y\}$ , where  
853  $H$  is the centering matrix,  $K_Y$  is a kernel matrix of  $Y$ , and  
854  $K_{ZZ}$  is a kernel matrix of the mapped data. When using a

linear kernel, it becomes  $K_{ZZ} = Z^T \Gamma \Gamma^T Z$ . Therefore,  $\Gamma$   
can be calculated by solving the following optimization

$$\arg \max_{\Gamma} \text{Tr}\{\Gamma^T Z H K_{YY} H Z^T \Gamma\}, \quad (12)$$

$$\text{subject to } \Gamma^T \Gamma = I \quad (13)$$

This optimization has a closed-form solution, and the  
columns of the optimal  $\Gamma$  are the eigenvectors of  
 $M := Z H K_{YY} H Z^T$  corresponding to the  $d$  largest eigen-  
values where  $d$  is the dimensionality of the subspace [35].  
Here, since we only need a direction vector, we choose the  
eigenvector  $\hat{v}_1$  associated with the largest eigenvalue of  $M$ .

$$\delta A = \Gamma = \hat{v}_1. \quad (14)$$

To capture the non-linear relations of the attribute, a non-  
linear kernel can be used to calculate  $K_Y$  and  $K_{ZZ}$ . The  
closed-form solution for the non-linear case is provided in  
§ A.6.2.

### 870 A.6.2. Non-Linear $\delta A$

871 If we are interested in finding non-linear relations between  
872 the images in order to find direction of change in an attribute,  
873 a non-linear version of the formulation mentioned in the pre-  
874 vious subsection can be used. In this approach, similar to the  
875 linear case, we generate two sets of images associated with  
876 positive ( $a^+$ ) and negative ( $a^-$ ) categories of  $A$ . The set of  
877 CLIP features of these images are denoted by  $Z_A = \{z_i\}_{i=1}^m$   
878 and their corresponding labels of  $A$  as  $Y_A = \{y_i\}_{i=1}^m$ . Then  
879 we find an orthogonal transformation matrix  $\Gamma$  that maps  
880 kernelized  $Z_A$  to a subspace that maximizes the variance of  
881 the labeled data using supervised principal component anal-  
882 ysis [5] that maximizes the dependency between the mapped  
883 data  $\Gamma^T K_{ZZ}$  and  $Y_A$ .

884 Hilbert-Schmidt Independence Criterion (HSIC) [26] is  
885 employed as the dependence metric where its empirical ver-  
886 sion is defined as  $\text{HSIC}^{\text{emp}} = \text{Tr}\{H K_{ZZ} H K_Y\}$ , where  $H$   
887 is the centering matrix,  $K_Y$  is a kernel of  $Y$ , and  $K_{ZZ}$  is the  
888 kernelized  $Z_A$  using a similarity measure of the mapped data.  
889  $\Gamma$  can be calculated by solving the following optimization

$$\arg \max_{\Gamma} \text{Tr}\{\Gamma^T K_{ZZ} H K_{YY} H K_{ZZ}^T \Gamma\}, \quad (15)$$

$$\text{subject to } \Gamma^T \Gamma = I \quad (15)$$

This optimization has a closed-form solution and  
the optimal solution for  $\Gamma$  are the eigenvectors of  
 $M := K_{ZZ} H K_{YY} H K_{ZZ}^T$  corresponding to the  $d$  largest  
eigenvalues where  $d$  is the dimensionality of the sub-  
space [35]. Here, since we only need a direction vector,  
we choose the eigenvector  $\hat{v}_1$  associated with the largest  
eigenvalue of  $M$ .

$$\delta A = \Gamma = \hat{v}_1. \quad (16)$$

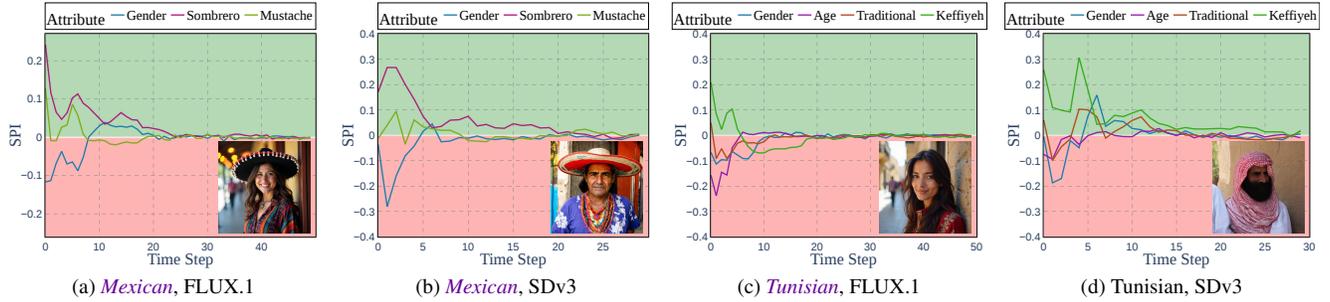


Figure 8. Additional samples of SPI plots of generated images by FLUX.1 and SDv3 for *Mexican* and *Tunisian* nationalities. A positive value for  $SPI(A, t)$  means that  $a^+$  is added to the image at time step  $t$ . Similarly, a negative SPI means that the image is moving toward  $a^-$ .

## 900 A.7. More Results on SPI

901 **More Sample-Wise Results.** More samples for SPI on  
902 *Mexican person* and *Tunisian person* are illustrated in Fig. 8.

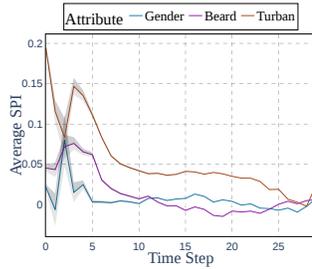


Figure 9. Average SPI in 100 samples for *Iranian person* generated by SDv3.

903 **Average SPI.** The average SPI in 100 images generated  
904 by SDv3 corresponding to “A photo of an Iranian person” is  
905 demonstrated in Fig. 9. As illustrated, the T2I model adds  
906 a high amount of information on attributes such as *turban*  
907 in the earlier steps. This confirms our earlier conclusions  
908 from sample-wise SPI in Fig. 6b. Additionally, we note  
909 that the variance for SPI is small in the time step  $T = 0$ ,  
910 suggesting the stereotypical predispositions noted in § A.2.4.  
911 However, in the next few time steps, we see a slightly larger  
912 variance that indicates that these models tend to correct the  
913 stereotypical attributes added in the former time steps.

## 914 A.8. More Results on T2I models’ Stereotypical Pre- 915 dispositions

916 In this section, we provide additional results that show that  
917 T2I models are predisposed to create stereotypical images  
918 for various *nationalities*. In Fig. 10, we show additional  
919 results for FLUX.1 on *Iranian*, *Indian*, and *Mexican*  
920 nationalities. Similar to our observations in § A.2.4, we note that  
921 the images generated from the velocity at  $t = 0$  for *Iranian*  
922 person contain stereotypical attributes such as *beard* and *tur-  
923 ban*. Likewise, for images of *Indian* personality, we observe  
924 *vibrant clothing* (e.g., orange veil). In the images of *Mexican*

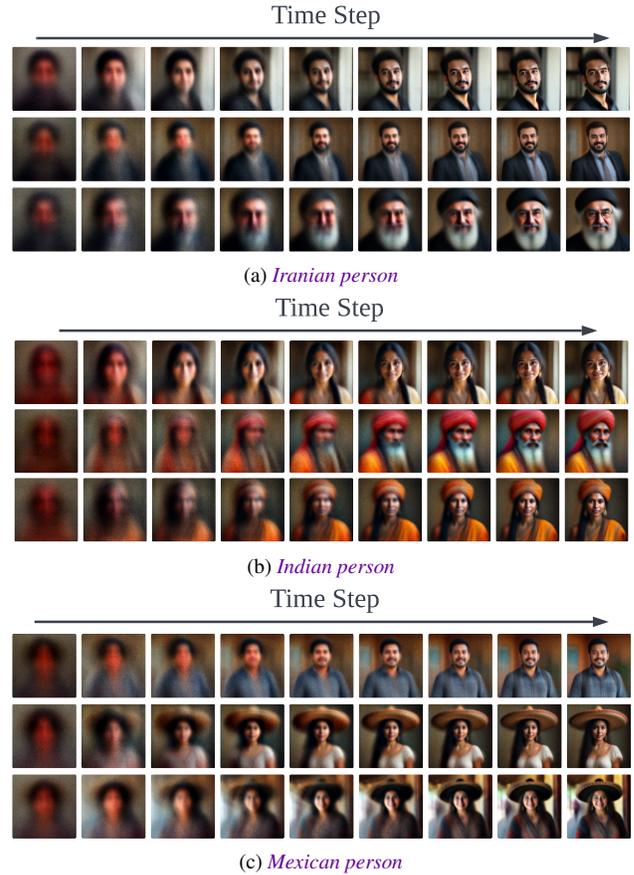


Figure 10. First 9 steps of image generation in FLUX.1 model for three nationalities: (a) *Iranian person*, (b) *Indian person*, and (c) *Mexican person*

person, we can see a faded *sombrero* in the early images. 925  
Moreover, the attributes that appear in the early stages of 926  
image generation are absent in the final generated image, 927  
indicating that these stereotypical attributes arise due to their 928  
intrinsic association with the concept. 929

## 930 A.9. Limitations

931 **Obtaining  $P^*(A | C)$ .** Access to  $P^*(A | C)$  is a crucial  
932 component for any stereotype measuring method. As men-  
933 tioned in § A.5.3,  $P^*(A | C)$  is obtained from census data  
934 and online sources when they are available. We note that  
935 reliable sources may not be available for every attribute and  
936 changes in survey methods can affect the results. However,  
937 for most stereotype evaluation and mitigation applications,  
938 reliable data can be found from government and survey agen-  
939 cies.

940 **Use of CLIP.** As mentioned in Sec. 3.1, we obtain  $P(A |$   
941  $I_i, C)$  using attribute classifiers. Instead of training attribute-  
942 specific classifiers, a zero-shot predictor like CLIP [42] can  
943 be utilized. However, some attributes may be unfamiliar  
944 to the model, resulting in lower accuracy in detecting them  
945 within the images. Additionally, these models may be biased  
946 in terms of concepts such as ethnicity. With advancements  
947 in zero-shot prediction models and the introduction of more  
948 accurate versions, newer models can seamlessly replace the  
949 existing ones in OASIS, thanks to its modular design.

950 For a small dataset of *doctors* that is used in Tab. 2, we  
951 evaluate the performance of the CLIP model in predicting  
952 the *gender*. As mentioned earlier, for each T2I model, we  
953 generated 100 images of *doctors* and manually labeled their  
954 genders. The accuracy of the CLIP model in predicting  
955 *gender* is demonstrated in Tab. 4. The results suggest that  
956 on this small dataset, the CLIP model can predict *gender*  
957 almost as accurately as human annotators.

Table 4. Performance of the CLIP model on predicting *gender* in the generated *doctors* dataset.

Model	Accuracy
SDv2	100%
SDv3	99%
FLUX.1	99%

958 Additionally, we evaluate the performance of the em-  
959 ployed CLIP model on CelebA [32] that contains more than  
960 200,000 face images of celebrities annotated with 40 binary  
961 attributes. Since we primarily used CLIP to predict attributes  
962 such as *beard* and *hat*, we evaluate the model on similar  
963 attributes (i.e., *having beard*, *man*, *wearing a hat*, *having*  
964 *a mustache*). The accuracy in predicting each attribute is  
965 reported in Tab. 5. The results demonstrate that the CLIP  
966 model can predict the attribute with an acceptable accuracy.  
967 As we noted earlier, although the CLIP model may not accu-  
968 rately predict certain general attributes, our results indicate  
969 that the CLIP model is suitable for predicting the attributes  
970 that we considered in this work.

971 The above-mentioned experiments, show the effective-  
972 ness of using a CLIP model in automating the classification

Table 5. Performance of the CLIP model on predicting four at-  
tributes in CelebA dataset.

Attribute	Accuracy
<i>having beard</i>	83.06%
<i>gender</i>	99.38%
<i>wearing a hat</i>	96.14%
<i>having a mustache</i>	94.77%

of the images. However, the accuracy of the model is not  
100% which indicates that by newer vision-language model  
with higher accuracy compared to the CLIP model that is  
employed in this paper, should be replaced in the OASIS.

## A.10. Importance of Detecting Stereotypes in Gen- erative Models

Visual content produced by generative models inadvertently  
perpetuates stereotypes about various ethnicities, cultures,  
nationalities, and professions [3]. Such images and videos  
are shared on online social media accounts such as X and  
Reddit, and this can reinforce stereotypical notions about  
certain social groups. This content could also influence  
public perception of marginalized communities and could  
undermine ongoing efforts to integrate them into mainstream  
society. For example, it has been noted that generated images  
of women from certain ethnicities tend to be sexualized [30].  
Additionally, the adoption of these generative models by  
various companies and institutions may have unforeseen  
consequences. For example, Nicoletti and Bass [37] states  
that using generative AI to develop suspect sketches could  
lead to wrongful convictions.

## A.11. Qualitative Descriptions of the Generated Datasets

We evaluated OASIS on the images corresponding to vari-  
ous nationalities generated by different T2I models. In this  
section, we give a qualitative description of the generated  
images. A few *randomly* selected representative samples  
from each culture and T2I model are shown in Fig. 11.

**FLUX.1.** The images produced by FLUX.1 are of high  
quality and look realistic. However, some stereotypes can be  
qualitatively observed from the images in Fig. 11a. For exam-  
ple, some images of *American* people contain the *American*  
*flag*. Images of *Indian* people tend to show vibrant colored  
clothing. Most of the generated images of *Iranian* people are  
of *men* and most of them wear *turban*. Among the images  
of *Mexican* people, *sombrero* is the most common stereo-  
typical element. We additionally note a general superficial  
diversity among the samples. For example, the images in  
each row were generated with the same random seeds. We  
can observe that the backgrounds in these photos are some-  
times repeated. For instance, compare the first columns of

1014 *Indian* and *Mexican* samples. Additionally, there is a clear  
 1015 disparity in the backgrounds across nationalities. The back-  
 1016 grounds for *American* and *Iranian* images are more often  
 1017 indoors than for *Indian* and *Mexican*. Images of *American*  
 1018 and *Iranian* people also more often contain images of offi-  
 1019 cials compared to *Indian* and *Mexican* images. Interestingly,  
 1020 Donald Trump’s image appeared when prompted to generate  
 1021 images of *American* person.

1022 **SDv3.** Fig. 11b shows the samples generated by SDv3.  
 1023 Although the generated images are of high fidelity, unlike  
 1024 FLUX.1, they lack variety in background and poses. Sur-  
 1025 prisingly, images of *American* people are relatively free of  
 1026 stereotypes and show ethnic diversity. However, the face  
 1027 images of *Indian* people look very similar and contain ele-  
 1028 ments such as *tilak/bindi*. The diversity drops further in the  
 1029 images of *Iranian* people. All the randomly selected samples  
 1030 contained images of *men* wearing *turban* and *religious at-*  
 1031 *tire*. Among the images of *Mexican* people, *sombreros* were  
 1032 present but in fewer proportions compared to the images  
 1033 generated by FLUX.1.

1034 **SDv2.** Some representative samples generated by SDv2  
 1035 are shown in Fig. 11c. Among the considered T2I models,  
 1036 SDv2 produced images with the least photorealism, with  
 1037 some displaying distorted facial expressions. However, these  
 1038 images generally contain diverse facial attributes such as  
 1039 hairstyle. Images of *American* people are of higher qual-  
 1040 ity compared to other nationalities, although they include  
 1041 black & white portraits. We note the lack of ethnic diver-  
 1042 sity among these images compared to those in SDv3 and  
 1043 FLUX.1. Although identity diversity is lower for images  
 1044 of *Indian* people compared to FLUX.1 and SDv3, we also  
 1045 observe fewer stereotypical attributes. Similar to SDv3, the  
 1046 images of *Iranian* people generated by SDv2 are primarily of  
 1047 *men*, mostly donning *turban*. Stereotypes such as *sombrero*  
 1048 and *colorful clothing* are present in the images of *Mexican*  
 1049 people. Among all the T2I models that we considered, SDv2  
 1050 seems to have the least gender diversity across all national-  
 1051 ities.

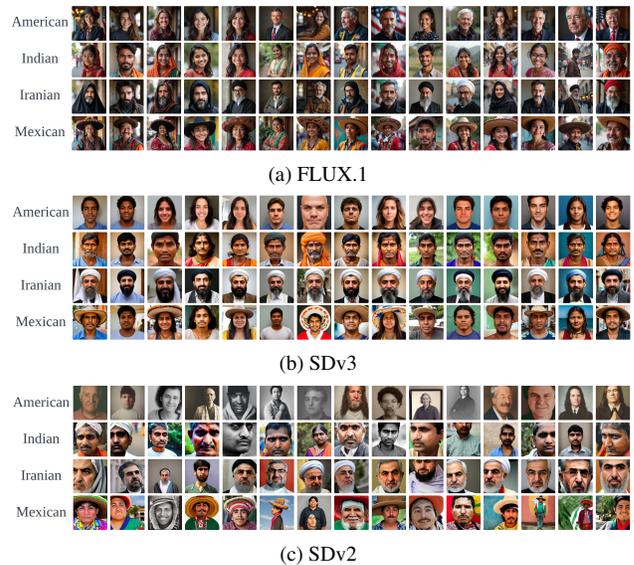


Figure 11. A few *randomly* selected representative samples from each culture and T2I model.