

GradBias: Unveiling Word Influence on Bias in Text-to-Image Generative Models

Moreno D'Incà, Elia Peruzzo, Massimiliano Mancini, Xingqian Xu, Humphrey Shi, Nicu Sebe

Abstract—Recent progress in Text-to-Image (T2I) generative models has enabled high-quality image generation. As performance and accessibility increase, these models are gaining significant attraction and popularity: ensuring their fairness and safety is a priority to prevent the dissemination and perpetuation of biases. However, existing studies in bias detection focus on closed sets of predefined biases (e.g., gender, ethnicity). In this paper, we propose a general framework to identify, quantify, and explain biases in an open set setting, *i.e.* without requiring a predefined set. This pipeline leverages a Large Language Model (LLM) to propose biases starting from a set of captions. Next, these captions are used by the target generative model for generating a set of images. Finally, Vision Question Answering (VQA) is leveraged for bias evaluation. We show two variations of this framework: OpenBias and GradBias. OpenBias detects and quantifies biases, while GradBias determines the contribution of individual prompt words on biases. OpenBias effectively detects both well-known and novel biases related to people, objects, and animals and highly aligns with existing closed-set bias detection methods and human judgment. GradBias shows that neutral words can significantly influence biases and it outperforms several baselines, including state-of-the-art foundation models. Code available here: <https://github.com/Moreno98/GradBias>.

Index Terms—Text-to-Image Generation, Fairness, Bias

1 INTRODUCTION

Text-to-Image (T2I) generation has caught the attention of the general public thanks to the recent advancements in quality and fidelity, reaching photo-realistic results [1]–[5]. These advancements unlocked multiple use cases, including image editing [6]–[9], personalization [10], [11] and additional conditioning [12]–[14] allowing even unexperienced users to obtain high-quality images. With this growing popularity, it is becoming increasingly important to investigate the safety and fairness of such models to avoid unwanted and unexpected behaviours [15]–[17]. Studying ethical issues in T2I generative models is an open area of research as defining the concept of *bias* presents significant challenges due to its inherently subjective nature and diversity [16]. As biases can range from more visual ones (e.g., gender and ethnicity), to more complex ones (e.g., socioeconomic status and cultural representation), establishing a universal definition remains challenging [18].

Given a concept t and a set of classes \mathcal{C} describing t , we define a T2I model G as biased against the concept t if the generated images do not follow a uniform distribution across \mathcal{C} [16], [20]. Therefore G tends to favor a particular class $c \in \mathcal{C}$ (e.g., “man”) when provided with a class-neutral textual prompt (e.g., “A chef cooking.”). Over the years, bias mitigation techniques have been proposed to control or remove specific biases without compromising performance [15], [16], ranging from training-time [21]–[24] to data augmentation methods [25], [26]. Nevertheless, applying these techniques poses a significant challenge, as retraining T2I models is often prohibitively resource-intensive.

Consequently, inference-time methods have gained attention, including fairer guidance methods [16], [27], [28] and prompt learning [15]. While these methods yield great results, they all require a pre-defined set of biases, limiting their applicability to well-known concepts such as gender, age, ethnicity [16], [17], [27], [28] and facial attributes [15].

We argue that T2I models may exhibit unknown biases not captured by close-set detection pipelines. The example in fig. 1 illustrates an overview of the close-set and open-set frameworks. A caption “A person using a laptop in the kitchen” allows close-set pipelines to identify predefined biases such as gender, age, and race, but may miss other significant biases like laptop brand or kitchen style. This limitation prompts the question: “Can we detect biases in an open-set fashion?”. This novel task poses challenges, including the difficulty of identifying all potential biases and the prohibitive cost of data collection.

We introduce a general pipeline for discovering and evaluating biases in T2I models, without requiring predefined categories. This pipeline exploits a Large-Language-Model (LLM) to propose biases starting from textual captions. Afterward, it queries a Visual Question Answering (VQA) model with the generated images and the potential biases, previously identified. This design removes the need for attribute-specific classifiers, which are unsuitable for open-set settings, thereby addressing another key limitation of previous works [15], [16], [27]–[29].

We implement two variations of this pipeline: OpenBias and GradBias. OpenBias studies the general behavior of T2I models, querying the VQA to evaluate the presence of the proposed biases. GradBias analyzes biases at the sentence level, effectively modeling how different words in a sentence contribute to the overall bias. This reflects a practical scenario, mirroring the standard use case of such models, where users seek to control the generation

M. D'Incà, E. Peruzzo, M. Mancini, and N. Sebe are with the University of Trento (Trento, Italy). X. Xu and H. Shi are with SHI Labs @ Georgia Tech & UIUC & Picsart AI Research (PAIR). Corresponding authors: H. Shi (shi@gatech.edu) and N. Sebe (niculae.sebe@unitn.it)

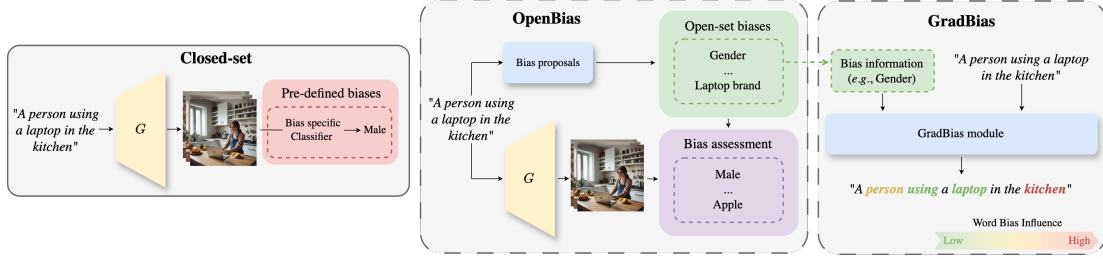


Fig. 1: We propose a modular framework for open-set bias evaluation. In contrast to previous works [15], [16], [19], our pipeline does not require a predefined list of concepts but proposes a set of novel biases. We implement two variations: OpenBias discovers general biases in T2I models, while GradBias detects the influence of each prompt word on the bias.

based on a specific prompt. For example, a gender-neutral word like “kitchen” might influence the model to produce a gender-biased output (as shown in fig. 1). Investigating these spurious correlations can offer deeper insights into T2I models, aiding researchers in creating fairer systems.

The contributions of this work are:

- To the best of our knowledge, we are the first to address the problem of bias detection in an open-set fashion without requiring a predefined list of biases.
- We present a novel framework specifically designed for open-set bias evaluation. This pipeline leverages an LLM to build a knowledge base of biases and a VQA model for bias assessment.
- We propose two variations of this pipeline: OpenBias and GradBias. OpenBias studies biases at high-level, automatically detecting well-known biases and uncovering novel ones. GradBias identifies the learned spurious correlations between individual words and the generated content, providing insights into the role of specific words in contributing to the bias.
- We evaluate both OpenBias and GradBias on multiple versions of Stable Diffusion [16], [27], [30]. OpenBias demonstrates its agreement with closed-set bias detection methods and human judgment, while GradBias outperforms several newly introduced baselines including those based on large multi-modal models.

This paper extends our previous work [31] by drawing a link between bias discovery and explanation under the same, unified framework. In this regard, we introduce GradBias, a novel method specifically tailored to determine the influence of individual prompt words on biases in T2I models. Specifically, while assessing the bias with the VQA model, GradBias computes the cross-entropy loss between the VQA’s output and the attribute representing the targeted bias. Through standard backpropagation of this loss, we can estimate the relative contribution of each token of the conditioning prompt with respect to the bias under study. As this research direction is under-explored, we introduce several baselines, exploiting recent advancements in LLMs and large multi-modal models. GradBias outperforms such baselines and unveils the important correlation between neutral words and bias in T2I models. Beyond introducing GradBias, we also expand the related work section to include works modeling language influence on bias and revise the method section to make the link between OpenBias and GradBias explicit.

2 RELATED WORK

Pipeline with Foundation Models. Large-scale deep learning models, often referred to as foundation models [32], are trained on extensive datasets and have demonstrated remarkably high performance in zero-shot settings [33]–[35]. These models are typically trained with self-supervision objectives [32] making them suitable across different modalities, including text [33], [36], vision [37]–[39] and even multiple modalities simultaneously [40]–[42]. These advancements enabled us to achieve complex tasks previously unthinkable. Recent works have demonstrated the capabilities of LLMs to generate Python code for querying vision and language models [43], [44], enabling the creation of automatic evaluation pipelines. For instance, TIFA [45] evaluates the image-text faithfulness of T2I models by querying a VQA model with LLM-generated questions. Moreover, captioning can be improved by querying a VQA model with LLM-generated questions, leading to more accurate and detailed descriptions [40], [46]. A similar study [47] identifies spurious correlations in generated images through captioning and model interpretation. However, the discovered biases are not quantified or categorized.

Similarly to the above pipelines, our general framework exploits foundation models for automatic open-set bias detection. It constructs a knowledge base of biases using an LLM and assesses them via VQA.

Bias detection in generative models. Bias mitigation in generative models has been extensively studied, particularly with GAN-based approaches. [48] alters the latent space semantic distribution at inference time to improve fairness. [19] ensures fairer representations for sensitive groups with gradient clipping. The research community has recently shifted its attention towards T2I models. In FairDiffusion [16] user-provided instructions are used to add a fair guidance term to enhance classifier-free guidance [49], effectively guiding Stable Diffusion [4] toward fairer generation in job-related domains. Similarly, *ITI-GEN* [15] exploits prompt learning to make T2I models fairer in the context of handwritten digits. In [29] unsupervised learning estimates the data manifold of the training set to guide the generation towards fairer representations. Recent works mitigate biases by aligning generated images with a user-defined target distribution. For example, [27] uses distribution guidance, while [28] employs distributional alignment loss and finetuning to achieve this alignment. Another line of research focuses on mitigating inappropriateness in T2I generation. [50] shows how (negative) prompt and semantic guidance [51] can be used to mitigate such behaviour.

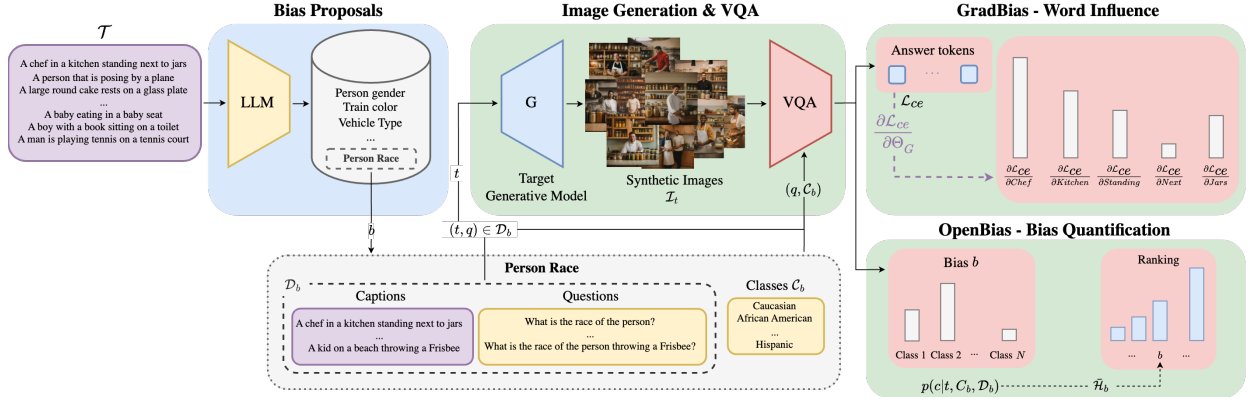


Fig. 2: We propose a general pipeline for open-set bias detection, quantification, and word-level explanation. Starting with a dataset of real textual captions (\mathcal{T}) we use a Large Language Model (LLM) to build a knowledge base \mathcal{B} of potential biases occurring in the image generation process. Next, the target generative model synthesizes images using captions where a potential bias has been identified. Finally, Vision Question Answering (VQA) is employed for either bias quantification or word-level bias explanation, depending on whether OpenBias or GradBias is applied.

These bias mitigation methods share a notable limitation: defining biases beforehand, limiting their applicability to well-known concepts. We argue that unconsidered and unstudied biases may still exist in T2I models. *Therefore, the framework introduced in this work serves as a complementary tool, enhancing the utility of existing methods by identifying, quantifying, and improving the understanding of additional biases.*

Language influence on biases. While existing research has explored biases in T2I models, there remains a significant gap in detecting how individual prompt words influence these biases. The first work in this direction is [52] which assesses the influence of each word by iteratively replacing it with synonyms and recalculating bias metrics. While this method may exhibit high performance, it demands high computational resources as it requires regenerating the set of images each time a word is replaced. We refer the reader to 4.2.1 where we treat this approach as upper-bound and use it for determining the ground truths in our experiments. Another approach [53] achieves targeted demographic distributions by using an LLM to modify input prompts. However, while effective in guiding demographics, it does not assess the impact of individual prompt words on bias. [54] improves output controllability by manipulating the language embedding space, increasing the expressivity beyond standard language descriptions. This approach generates images that match specific attribute distributions, however, it does not assess the impact of individual prompt words on the bias.

These works improve model explainability and controllability. However, a *direct* method to evaluate the impact of individual prompt words on bias remains unexplored. Hence, we propose a variation of our general pipeline to *directly* assess word-level bias influence and introduce baselines for benchmarking this task.

3 OPENBIAS AND GRADBIAS

OpenBias identifies and quantifies general biases in T2I models, while GradBias analyzes the impact of individual prompt words on these biases. Both methods share a common pipeline, illustrated in Figure 2. The bias proposal stage queries an LLM with real textual captions to build a

knowledge base of potential domain-specific biases, eliminating the need for predefined concepts and allowing the discovery of new biases. These captions are then used for generating images with the target generative model. Finally, a Vision-Question-Answering (VQA) model evaluates the potential biases. We describe the general framework and its two variations below.

3.1 Bias proposal module

The bias proposal module aims at building a knowledge base \mathcal{B} composed of domain-specific biases. This stage queries an LLM with real captions \mathcal{T} to produce: the potential bias name, the set of classes associated with that bias, and a question for bias assessment.

Formally, for a given caption $t \in \mathcal{T}$, we represent the LLM’s output as a set of triplets $\mathcal{L}_t = \{(b_i^t, \mathcal{C}_i^t, q_i^t)\}_{i=1}^{n_t}$, where the cardinality of the set n_t is caption-dependent. This triplet (b, \mathcal{C}, q) is composed of the proposed bias b , a set of associated classes \mathcal{C} , and the question q specific to the caption t . We follow the in-context learning paradigm [33], [55] and prompt the LLM with task descriptions and examples (see the Appendix for the exact system prompt). We aggregate the above LLM output for the whole dataset by defining the caption-specific biases B_t as the union of its potential biases *i.e.* $B_t = \bigcup_{i=1}^{n_t} b_i$. Therefore, the knowledge base of biases \mathcal{B} is built as the union of the caption-level ones, *i.e.* $\mathcal{B} = \bigcup_{t \in \mathcal{T}} B_t$. Afterward, we define a database \mathcal{D}_b specific to the bias b by collecting bias-specific information, *i.e.* captions and questions:

$$\mathcal{D}_b = \{(t, q) \mid \forall t \in \mathcal{T}, (x, \mathcal{C}, q) \in \mathcal{L}_t, x = b\}. \quad (1)$$

Additionally, we define $\mathcal{T}_b = \{t \mid (t, q) \in \mathcal{D}_b\}$ as the collection of captions linked to the bias b , and \mathcal{C}_b as the union of the set of classes associated with the bias b within \mathcal{T} .

However, \mathcal{D}_b does not consider whether the classes are specified in the caption. For example, when generating “An image of a large dog”, the dog’s size should not be treated as one of the potential biases since it is already stated in the prompt. We address this by employing a two-stage filtering on \mathcal{D}_b . In the first phase, given a sample $(t, q) \in \mathcal{D}_b$, we ask the LLM whether the answer to the question q is explicitly

present in the caption t . In the second stage, we filter out captions that contain either a class $c \in \mathcal{C}_b$ or its synonyms. We identify the synonyms for \mathcal{C}_b using ConceptNet [56]. We find that combining these two stages empirically produces more robust results. This module generates bias proposals in an open-set fashion tailored to the given dataset. OpenBias and GradBias will evaluate these biases at high-level and word-level respectively.

3.2 Image Generation & VQA

After building the knowledge base of potential biases \mathcal{B} , we evaluate them using the target generative model G . For a given bias $b \in \mathcal{B}$ and caption $t \in \mathcal{T}_b$, we generate a set of N images \mathcal{I}_b^t , defined as:

$$\mathcal{I}_b^t = \{G(t, s) | \forall s \in S\} \quad (2)$$

where S is the set of sampled random noise vectors with cardinality $|S| = N$. By sampling multiple noise vectors, we obtain a distribution of G 's output on the same prompt t .

In the final stage, the pipeline evaluates the bias b in \mathcal{I}_b^t by querying a state-of-the-art Vision-Question-Answering VQA model. For each image $I \in \mathcal{I}_b^t$, the VQA model answers the question q associated with prompt t from the pair $(t, q) \in \mathcal{D}_b$, selecting the answer from the available classes \mathcal{C}_b . Hence, we define the predicted class \hat{c} as:

$$\hat{c} = \text{VQA}(I, q, \mathcal{C}_b). \quad (3)$$

As shown in Figure 2, we implement two variations for this final step. OpenBias quantifies the overall biases identified in the generative model, while GradBias assigns scores to each input word, indicating its influence on the bias.

3.3 OpenBias - Bias Assessment and Quantification

In OpenBias, we quantify the severity of the proposed bias $b \in \mathcal{B}$ by applying the VQA process to the whole \mathcal{I}_b^t and draw a distribution over the classes \mathcal{C}_b . We then explore two different scenarios: *context-aware* where we examine bias in caption-specific images \mathcal{I}_b^t , and *context-free* where we consider the entire set of images \mathcal{I}_b associated with the bias.

3.3.1 Context-Aware Bias

In the context-aware formulation, we capture the bias \mathcal{B}_j by considering the caption context. As described in Section 3.1, the bias proposal module filters out captions that mention bias attributes, retaining only those neutral to the bias. However, other factors within the caption may influence the direction and intensity of the bias. For instance, the captions "A military is running" and "A person is running" are both neutral to the bias "person gender", but the bias's direction and intensity may differ based on the context.

While assessing the bias, we consider the role of the context by analyzing the set of images \mathcal{I}_b^t generated from a specific caption $t \in \mathcal{T}$, effectively collecting statistics at the caption-level. Consequently, given a bias b , the probability for a class $c \in \mathcal{C}_b$ is computed as:

$$p(c|t, \mathcal{C}_b, \mathcal{D}_b) = \frac{1}{|\mathcal{I}_b^t|} \sum_{I \in \mathcal{I}_b^t} \mathbb{1}(\hat{c} = c) \quad (4)$$

where $\hat{c} = \text{VQA}(I, q, \mathcal{C}_b)$ is the prediction of the VQA as defined in eq. (3), and $\mathbb{1}(\cdot)$ is the indicator function.

3.3.2 Context-Free Bias

Here, we are interested in studying the overall behavior of the model G , offering insights into aspects such as the majority class (*i.e.* the direction of the bias) and the overall bias intensity. We exclude the caption context by averaging the VQA scores for $c \in \mathcal{C}_b$ over all captions t related to that bias $b \in \mathcal{B}$:

$$p(c|\mathcal{C}_b, \mathcal{D}_b) = \frac{1}{|\mathcal{D}_b|} \sum_{(t,q) \in \mathcal{D}_b} p(c|t, \mathcal{C}_b, \mathcal{D}_b) \quad (5)$$

Note that the context-aware bias is a special case of this scenario, where \mathcal{D}_b contains a single instance, *i.e.* $\mathcal{D}_b = \{(t, q)\}$.

3.3.3 Bias Quantification and Ranking

The open-set nature of our setting challenges the study of the multiple biases. Therefore, we propose to rank them for a comprehensive analysis. Following previous works [16], [20], we consider the generative model G unbiased for a concept b if the distribution over the classes \mathcal{C}_b is uniform. We measure the bias intensity as the entropy of the class probability distribution, obtained with either eq. (4) or eq. (5). Since biases may have different cardinalities in \mathcal{C} , we normalize the entropy by its maximum value [57] to enable comparisons. Finally, we refine this score for readability as:

$$\bar{\mathcal{H}}_b = 1 + \frac{\sum_{c \in \mathcal{C}_b} p(c|\mathcal{C}_b, \mathcal{D}_b) \log p(c|\mathcal{C}_b, \mathcal{D}_b)}{\log(|\mathcal{C}_b|)} \quad (6)$$

where $\bar{\mathcal{H}}_b \in [0, 1]$, with 0 denoting low bias and 1 high bias.

3.4 GradBias - Word Influence

As discussed in Sections 1 and 3.3, bias in T2I models is affected by multiple factors, including caption context and specific words. We focus on specific words and introduce GradBias, a novel method for detecting the influence of prompt words on the bias. After applying the shared pipeline detailed in Sections 3.1 and 3.2 and shown in Figure 2, GradBias aims to identify how individual words influence the generation of specific attributes $y \in \mathcal{C}_b$ related to the proposed bias b . To achieve this, we first compute the following cross-entropy loss CE:

$$\mathcal{L} = \text{CE}(\text{VQA}(I, q, \mathcal{C}_b), y) \quad (7)$$

where $\text{VQA}(I, q, \mathcal{C}_b)$ denotes the logits from the VQA model over \mathcal{C}_b (note: this is an abuse of notation as in eq. (3) we have the same notation as argmax), and $y \in \mathcal{C}_b$ is the attribute related to the bias b . We select y as the argmax of the logits, as this attribute (i) describes the bias b and (ii) is present in the generated image I , according to the VQA's prediction. We then backpropagate this loss to compute the gradients of the text tokens conditioning the T2I model. Given a prompt t we denote the textual embeddings of its tokens as $\{e_0, e_1, \dots, e_n\}$. For each token embedding e_i , we compute the absolute value of the scalar gradients as:

$$e'_i = \left| \frac{\partial \mathcal{L}}{\partial e_i} \right| \quad (8)$$

Therefore, e'_i estimates the contribution of the token e_i on generating the specific attribute related to the bias b , according to eq. (7). The final score for a word $w \in t$ is the

Model	Gender		Age		Race	
	Acc.	F1	Acc.	F1	Acc.	F1
CLIP-L [42]	91.43	75.46	58.96	45.77	36.02	33.60
OFA-Large [58]	93.03	83.07	53.79	41.72	24.61	21.22
mPLUG-Large [59]	93.03	82.81	61.37	52.74	21.46	23.26
BLIP-Large [60]	92.23	82.18	48.61	31.29	36.22	35.52
Llava1.5-7B [41], [61]	92.03	82.33	66.54	62.16	55.71	42.80
Llava1.5-13B [41], [61]	92.83	83.21	72.27	70.00	55.91	44.33

TABLE 1: VQA evaluation on the generated images using COCO captions. In **gray** the chosen default VQA model.

gradients of the tokens representing that word, addressing cases where a single word is split into multiple sub-tokens. The intuition is that these gradient magnitudes reveal how changes in specific words affect the VQA model’s output and, consequently, they indicate the contribution of each word in influencing the generative model towards biased outcomes as defined by the classes C_b . For example, given the prompt “A chef in a kitchen” and the bias “gender”, the scores assigned to each word estimate their contribution in generating a specific bias-related attribute “male” or “female”. We show a result related to this example in Figure 16.

4 EXPERIMENTS

We assess OpenBias by (i) comparing it against a state-of-the-art closed-set bias detection classifier and (ii) evaluating its alignment with human judgment via a user study. GradBias is evaluated against several baselines and state-of-the-art large vision and language models.

4.1 OpenBias

In OpenBias, we heavily rely on the zero-shot capabilities of each module. We follow [16] and assess its performance by comparing it against FairFace [62], a state-of-the-art classifier on gender, age, and race. Moreover, we evaluate the alignment of OpenBias with human judgment via a user study. Below, we define the implementation details and show the results of these experiments.

4.1.1 Pipeline implementation

Datasets. We use captions from two multimodal datasets: Flickr 30K [63] and COCO [64]. Flickr 30K [63] comprises 30k in-the-wild images with 5 captions per image. Similarly, COCO [64] is a comprehensive dataset spanning a diverse range of images with complex contexts. We filter COCO by selecting captions referring to a single person only, leading to roughly 123K captions. This extensive subset targets bias detection in the context of people, a critical area for examining biases. Nevertheless, it is worth noting that OpenBias also identifies biases outside this domain to include objects, animals, and actions, as shown in Section 5.1.

Implementation details. The flexibility and modularity of OpenBias allow the replacement of individual modules as soon as better models are available. In this study, we use Llama2-7B as our LLM of the bias proposal module and Llava1.5-13B as our Vision-Question-Answering (VQA) model. This choice follows our evaluation described in section 4.1.2, where Llava1.5-13B is the best-performing model. While running OpenBias, we randomly sample 100 captions associated with each bias and generate 10 images for each caption, leading to 1000 generated images for each bias.

Model	Flickr 30k [63]			COCO [64]		
	gender	age	race	gender	age	race
Real	0	0.032	0.030	0	0.041	0.028
SD-1.5 [4]	0.072	0.032	0.052	0.075	0.028	0.092
SD-2 [4]	0.036	0.069	0.047	0.060	0.045	0.105
SD-XL [5]	0.006	0.028	0.180	0.002	0.027	0.184

TABLE 2: KL divergence (\downarrow) computed over the predictions of Llava1.5-13B and FairFace on generated and real images.

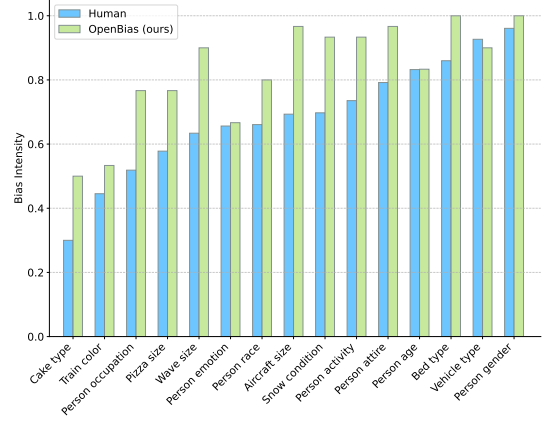


Fig. 3: Human evaluation results.

4.1.2 Quantitative results

Agreement with FairFace. We first evaluate the alignment of several SoTA VQA models with FairFace on generated images by Stable Diffusion XL [5]. We treat FairFace predictions as ground truths and report accuracy and f1-score in Table 1. Llava1.5-13B [41], [61] emerges as the best VQA model on gender, age, and race tasks, therefore, we employ it as OpenBias’s default VQA. We further test the robustness of Llava1.5-13B by applying the same evaluation method to both real and generated images. We use real pictures from Flickr 30K [63] and COCO [64] alongside images generated by Stable Diffusion 1.5, 2, and XL [4], [5]. The agreement between Llava1.5-13B and FairFace is measured using KL divergence between their respective probability distributions. The results, presented in Table 2, show low KL divergence, indicating strong alignment and confirming the VQA model’s robustness across real and generated images.

User Study. We evaluate OpenBias’s alignment with human judgment through a context-aware human study. We generate 10 for each caption and bias. We publish the survey on crowdsourcing platforms, without geographical restrictions, and randomize the questions’ order. Participants are asked to indicate the bias direction (majority class) and rate its intensity on a scale from 0 to 10, with an option for “No Bias” if none is perceived. This study comprises 15 diverse object-related and person-related biases, accounting for 390 total presented images. We receive 2200 valid responses from 55 unique individuals. We compare the participants’ intensity scores with OpenBias in Figure 3. We can observe high alignment on multiple biases, especially for “Person age”, “Person gender”, “Vehicle type”, “Person emotion”, and “Train color”. The study covers 15 object-related and person-related biases, resulting in 390 presented images. We receive 2200 valid responses from 55 unique individuals. We compare the participants’ intensity ratings with OpenBias in Figure 3, finding high alignment on multiple biases, including “Per-

Model	Stable Diffusion 1.5 [4]			Stable Diffusion 2 [4]			Stable Diffusion XL [5]		
	Top-1	Top-2	Top-3	Top-1	Top-2	Top-3	Top-1	Top-2	Top-3
Random	36.36	62.09	81.32	37.28	63.08	81.25	37.54	62.87	81.04
Dependency Tree	36.86	65.36	81.57	39.39	62.12	82.58	39.39	62.40	80.56
Llama2-13B [36]	35.14	64.37	79.85	41.67	64.39	79.04	41.18	66.50	80.05
Llama3-8B [36]	39.07	66.58	82.56	40.40	69.19	82.58	37.85	66.50	81.07
Llama3-70B [36]	39.07	63.64	81.57	39.14	67.42	83.59	41.18	65.98	81.33
Llava-1.5-13B [41], [61]	37.35	62.65	81.08	44.70	67.93	83.84	41.43	62.92	79.03
GradBias (CLIP-L)	48.65	71.50	86.98	47.98	78.79	89.14	44.25	66.24	84.65

TABLE 3: GradBias evaluation on Stable Diffusion 1.5, 2, and XL [4], [5].

son age”, “Person gender”, “Vehicle type”, “Person emotion”, and “Train color”. We also compute the Absolute Mean Error (AME) between the predicted model’s bias intensity and the average user score, resulting in an $AME = 0.15$. Additionally, we assess bias direction by comparing the agreement on the majority class, with OpenBias aligning with human choices 67% of the time. It is important to note that concepts of bias and fairness are highly subjective, potentially introducing errors in the evaluation process. Nevertheless, our results show a correlation between human and OpenBias scores, validating our pipeline.

4.2 GradBias

We evaluate GradBias against newly introduced baselines, ranging from natural language processing (NLP) techniques to methods leveraging advanced foundation models. Additionally, we perform an ablation study to examine various components of our pipeline. Below, we provide the implementation details and results regarding these experiments.

4.2.1 Pipeline implementation

Dataset. As discussed in section 1, understanding the impact of specific prompt words on the bias remains an under-explored area, resulting in a lack of task-specific data. To address this, we build a textual dataset to evaluate GradBias. This dataset includes 11 biases related to people and objects identified by OpenBias, with up to 50 randomly sampled captions per bias, leading to 391 captions.

Groundtruth computation. In GradBias, each word can influence bias to different extents, and multiple words may have similar impacts. For instance, in the prompt “a nurse in the hospital”, “nurse” may heavily influence the bias towards females, while “hospital” might have a neutral impact. On the contrary, in “a scientist in the laboratory”, both “scientist” and “laboratory” might equally reinforce a male bias. Additionally, this influence is model-dependent, varying with the generative model in use. These factors complicate the establishment of a definitive ground truth.

We compute a reliable ground truth by adapting the method from [52]. Their approach involves generating a distinct set of images for each word in the prompt by replacing it with synonyms, thereby assessing how each substitution influences the model’s bias. Given its robustness, we treat this method as an upper bound to compute our ground truths. However, generating a full set of images for every replacement is resource-intensive and impractical. To address this, we modify this method to enhance its efficiency. For a given prompt t and bias classes C_b , we first generate

10 images using the full prompt and establish a baseline distribution over C_b with OpenBias. We then iteratively remove one word at a time, regenerate the images, and observe changes in the class distribution. The influence score for each word is determined by summing the variations in class contributions. The ground truth is defined as the word whose removal causes the most significant shift in the class distribution compared to the full prompt. It is important to mention that we (i) account for scenarios where multiple words equally influence the bias by considering multiple words with equal top scores as the ground truth, and (ii) perform the ground truth computation for each tested model to account for model dependency. Additionally, we exclude (i) stop-words as they do not provide meaningful insights and (ii) words that directly describe or are synonymous with the bias, as we aim to study the influence of bias-neutral words. For example, when studying the bias “person gender” given the prompt “A person cooking in the kitchen.” we exclude the word “person”. This is because “person” is associated with “person gender”, and its removal would naturally result in a high influence score.

Implementation details. We experiment with Llava1.5-13B [41], [61] and CLIP [42] as VQA models for GradBias. In the version with CLIP, we produce an answer via cosine similarity between the generated images and bias-related classes C_b . To ensure a fair evaluation, we compute the ground truth using BLIP2 [60] as the VQA model within OpenBias. This approach removes the risk of introducing unintentional biases, avoiding potential circular reasoning associated with relying solely on Llava1.5-13B for both GradBias and ground truth computation. Finally, we apply GradBias at each denoising step of Stable Diffusion [4], [5] and average over the steps. This approach ensures capturing the word influence throughout the entire denoising process, as this influence may change at different stages. In Figure 4 we ablate on how different denoising steps affect GradBias.

4.2.2 Quantitative results

Baselines. We evaluate GradBias by introducing three baselines: (i) a random-based method, (ii) an NLP-based method, and (iii) two foundation-model-based methods. The random approach randomly selects a word from the prompt as the most influential and serves as a lower-bound baseline. The NLP method leverages spaCy to compute the sentence dependency tree, then rotates the tree around the subject. We then apply a breadth-first search (BFS) starting from this root, ranking words based on the order in which they are visited. We have empirically found that starting from the sentence subject improves performance. Finally, we intro-

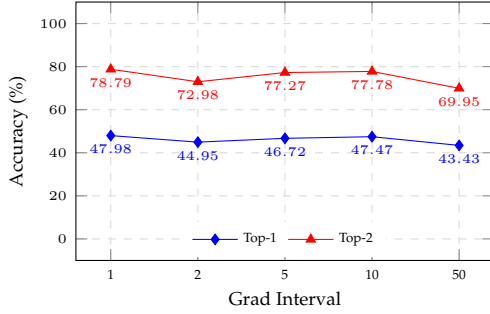


Fig. 4: Ablation on the gradient timing computation using Stable Diffusion 2 [4].

duce two foundation models based methods that perform answer ranking [33], [41], [65] on LLMs and VQAs models. For LLMs, we iteratively ask a yes or no question for each candidate word $t_i \in t$ to determine its influence on the bias b , ranking the words based on the “yes” logit probability. We structure the question as four variations of “Is **<candidate_word>** influencing **<bias>** bias in the prompt **<caption>**? Answer only with ‘yes’ or ‘no’.” and average over the answers’ probabilities. For VQA models, we follow the same procedure but leverage their multi-modal nature by providing the additional context of the generated images. We experiment with Llama2-13B, Llama3-8B, and Llama-70B [36] as LLMs and Llava1.5-13B [41], [61] as VQA.

Quantitative Results. We evaluate the baselines and GradBias with top-1, top-2, and top-3 accuracies across 11 biases, as described in Section 4.2.1. The results are presented in Table 3, where we apply the different methods to Stable Diffusion 1.5, 2, and XL [4], [5]. Across all settings, we observe a consistent progression from the random approach to GradBias. The dependency tree method outperforms the random approach but underperforms compared to other methods. This is expected as this approach is computed independently of the bias, making it less robust to bias diversity. In contrast, LLMs benefit from embedded bias knowledge showing a notable improvement. For example, Llama3-70B [36] outperforms the random baseline by +2.71% and +3.64% when applied to Stable Diffusion 1.5 and XL [4], [5], respectively. Similar gains are observed across other Llama configurations. Visual information provides additional insights that text-based approaches lack. Therefore, Llava-1.5-13B [41], [61] improves over the random baseline by +7.42% on Stable Diffusion 2. However, it aligns with LLMs in the other settings. In contrast, GradBias shows a significant improvement over the random approach, with gains of +12.29%, +10.70%, and +6.71% on Stable Diffusion 1.5, 2, and XL, respectively. It also outperforms all other baselines, achieving +9.58%, +3.28%, and +2.82% higher accuracy than the second-best method across all settings.

This evaluation underscores GradBias’s robustness against various baselines, including SoTA foundation models, demonstrating its effectiveness as a valuable tool for analyzing the impact of language on AI-generated content.

Ablation studies. As discussed in Section 4.2.1, evaluating the impact of computing the loss at different denoising steps is crucial. We assess this by analyzing images generated with Stable Diffusion 2 [4], computing the loss at intervals of 1, 2, 5, 10, and 50 denoising steps. As shown in Figure 4, GradBias performs optimally when the loss is computed

VQA	SD-2	SD-XL	SD-1.5
CLIP-L	47.98	44.25	48.65
Llava-1.5-13B	48.23	43.22	49.14

(a)

N	SD-2	SD-XL	SD-1.5
1	47.47	42.19	48.48
2	48.02	42.89	48.44
5	48.46	43.91	48.44
10	47.98	44.25	48.65

(b)

TABLE 4: We ablate on (a) the VQA used by GradBias, CLIP-L [42] or Llava-1.5-13B [41], [61] and (b) the number of images N considered by GradBias. We use images generated by Stable Diffusion 1.5, 2, and XL [4], [5].

multiple times during a single generation—*i.e.* at every 1, 2, 5, and 10 denoising steps. This is expected as GradBias accounts for the entire denoising process, capturing the complete semantic generation. In contrast, computing the loss only at the final 50-th step significantly reduces performance, with decreases of -4.55% in top-1 and -8.84% in top-2 accuracy. This highlights the importance of considering multiple denoising steps to achieve optimal results.

We further evaluate GradBias’s robustness by varying the VQA model. We test CLIP [42] and Llava1.5-13B [41], [61] and show the results in Table 4(a). Interestingly, GradBias performs consistently well regardless of the VQA used, suggesting that the choice of VQA has minimal impact on the computed gradients. This highlights that GradBias primarily focuses on the generative model, with the VQA serving primarily as a vehicle for gradient computation. This finding further validates our pipeline, as the gradients effectively capture the intrinsic biases of the generative model, minimizing the influence of the VQA.

In contrast to [52] where generating multiple images is *required* at each word replacement, GradBias should operate effectively with only one image since (i) it avoids computing distributions, and (ii) it relies on gradients to capture each word’s influence on the bias, independently of the number of generated images. We evaluate this by varying the number N of synthesized images, testing GradBias with N equals 1, 2, 5, and 10. Results in Table 4(b) show that GradBias performs optimally across all tested quantities, demonstrating its effectiveness even in low data regimes, proving advantageous in low-resource scenarios.

5 FINDINGS

In this section, we discuss the findings of OpenBias and GradBias when applied to three popular T2I generative models Stable Diffusion XL, 2, and 1.5 [4], [5].

5.1 OpenSet bias detection – OpenBias

We apply OpenBias to captions from COCO and Flickr30k, as detailed in Section 4.1, to analyze various biases across models and compare context-free versus context-aware bias.

Rankings. In Figure 5, we present the rankings of the multiple biases identified by OpenBias. Our method unveils both well-known biases (*i.e.* , “person gender”, “person race”) and novel ones (*i.e.* , “cake type”, “bed type”, “laptop brand”). By comparing the generative models we note a correlation in the detected biases, with the XL version amplifying them more than its predecessors. Moreover, OpenBias uncovers

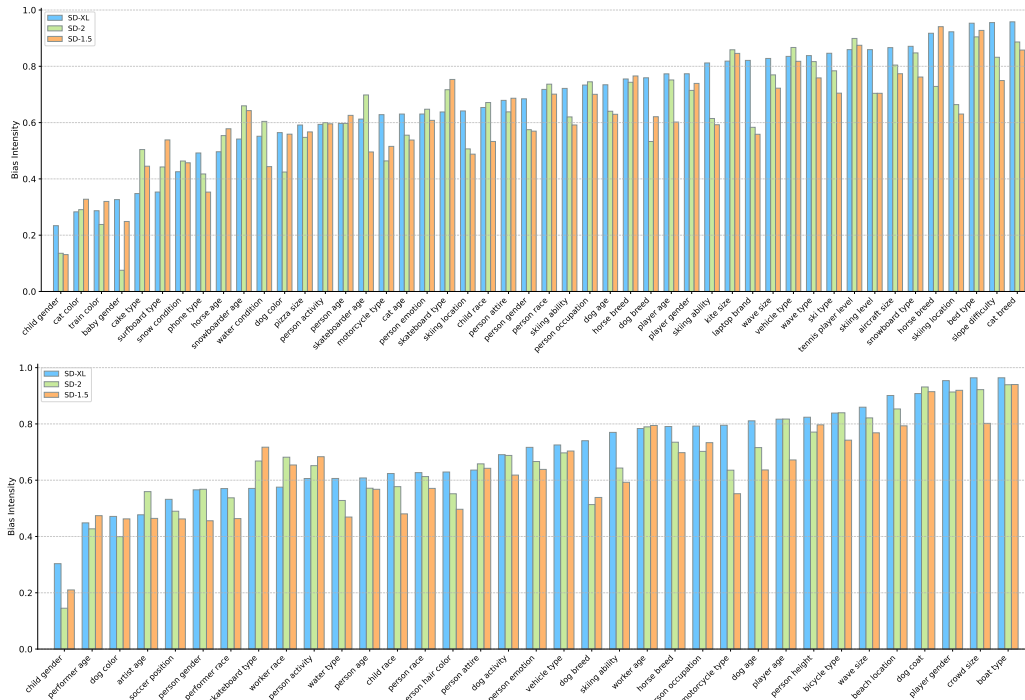


Fig. 5: Context-aware discovered biases on SD XL, 2 and 1.5 [4], [5] with COCO [64] and Flick30k [63] captions respectively.

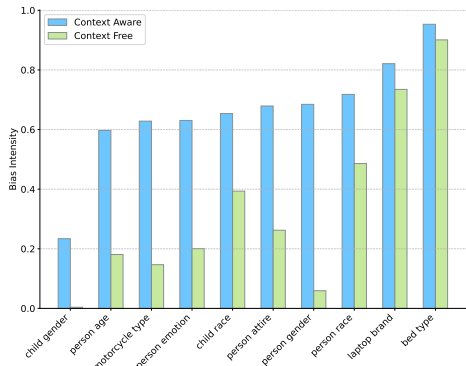


Fig. 6: Highlighting the importance of the context on Stable Diffusion XL [5] with the captions from COCO.

more object-centric biases from Flickr captions and person-centric biases from COCO. This diversity highlights OpenBias’s ability to propose domain-specific biases based on the dataset, demonstrating its adaptability to various domains.

Context-Free vs Context-Aware. We evaluate how additional elements in the caption affect biases by comparing the generative models in context-free and context-aware settings (see Section 3.3). In Figure 6, we compare these settings using Stable Diffusion XL. The low correlation between scores highlights the importance of context. For example, the bias “*motorcycle type*” is significantly higher in the context-aware setting, suggesting that the model generates specific motorcycle types based on the context. Differently, the “*bed type*” bias exhibits correlated intensity scores, indicating that the model generates the same bed type regardless of context.

Qualitative results. In Figure 7, we provide examples of biases identified by OpenBias in Stable Diffusion XL. We consider the context-aware setting, displaying images generated from the same caption. We highlight previously unexplored biases related to objects and animals, novel biases

concerning people, and well-known social biases. In the first row, the model consistently generates “yellow” trains and “quarter horses” from neutral captions. Additionally, it shows a pronounced bias towards the “Apple” brand, often generating laptops with a distinct “Apple” logo. OpenBias can identify novel person-related biases as Stable Diffusion XL frequently generates people in formal attire, demonstrating a “person attire” bias even with a neutral prompt. It is interesting to show that well-known social biases are also reflected in children. For instance, in the second row, the generative model associates a black child with an economically disadvantaged setting, described in the caption as “a dirt road”. This correlation between racial identity and socioeconomic status reinforces harmful stereotypes, highlighting the necessity of addressing novel biases in bias mitigation frameworks. Finally, OpenBias can detect well-known biases, such as “person gender” and “race”. Stable Diffusion XL exclusively generates “male” officers, even when provided with a gender-neutral prompt. We observe a “race” bias, with the generative model depicting only black individuals for “a man riding an elephant”. We provide further qualitative results in the Appendix.

Examining biases in the context-aware setting allows for a thorough investigation of emerging novel and socially sensitive biases. Our findings highlight the need for more inclusive open-set bias detection frameworks that can adapt to a wider range of scenarios, potentially contributing to the development of fairer generative models.

5.2 How does the bias arise? – GradBias

We qualitatively evaluate GradBias by applying it to Stable Diffusion XL [5] on multiple biases and present the results in Figure 16. We assess the influence of the top-scoring word identified by GradBias (highlighted in **red**) by manually replacing it with a different word (highlighted in **blue**),



Fig. 7: Biases discovered on Stable Diffusion XL [5] by OpenBias.

making sure to keep the overall sentence bias-neutral. In the first example, we consider the sentence “A chef in a kitchen standing next to a counter with jars and containers” and replace the GradBias predicted word “chef” with the more general word “person”. This replacement significantly shifts the bias from a male to a female representation, confirming that “chef” conveys a strong male representation. Similarly, replacing the predicted word “race” with “broken” in the sentence “A guy riding a race bicycle making a turn” results in individuals with darker skin tones. This shift highlights how associations between economic status and racial identity can reinforce and perpetuate harmful biases. GradBias also reveals the significant impact of context-specific words on bias. For example, replacing “computer” with “radio” in “A man uses his computer while sitting at a desk” shifts the generation toward older people. Similarly, replacing “book” with “comic” in “A man sitting by himself reading a book on a bench” results in more casual attire.

We build on OpenBias findings and apply GradBias to explore “child race” and “child gender” biases revealing that individual words significantly influence these biases. For example, in the sentence “A child is playing on a tennis court” replacing “tennis” with “basketball” shifts the generation towards a darker skin tone, showing the presence of sports-related stereotypes. Similarly, “child gender” bias is strongly affected by adjectives. In the sentence “A young child who is eating some food”, replacing “young” with “beautiful” shifts the gender representation from male to female, highlighting how appearance-related adjectives can significantly influence gender representations.

Finally, we apply GradBias to object-related biases, focusing specifically on “train color” and “laptop brand” biases. In the case of “train color”, Stable Diffusion XL generates trains with diverse colors when prompted with “A passenger train on a track next to a station”. However, when the GradBias predicted word “passenger” is replaced with “cargo”, the model predominantly generates yellow trains, demonstrating a shift toward a single color. Similarly, in “A photo of a person on a laptop in a coffee shop”, replacing the word “person” with “teenager” leads to more depictions of MacBook-like laptops, highlighting a brand bias associated with age-related stereotypes.

These results highlight the importance of studying the impact of individual words in shaping biases, as seemingly minor changes in the prompt can significantly influence biases. Moreover, word choices can introduce unwanted biases in the generation. It is important to note that multiple words in the same prompt can affect the bias differently. In this evaluation, we focus on replacing the top-scoring words from GradBias, however, other words can have similar impacts. For example, the word “kitchen” in “A chef in a kitchen standing next to a counter with jars and containers” is a strong candidate for shaping gender bias, as demonstrated by the counterpart prompt with “person” as the subject. Nevertheless, GradBias covers this scenario by providing a ranking based on the influence of each word on the bias, thus providing a global picture of what words should be considered meaningful for a given bias. Understanding the linguistic influence on biases is crucial for developing more fair and accurate generative models.



Fig. 16: Qualitative results of GradBias. We highlight in **red** the most influential word identified by GradBias.

6 LIMITATIONS

OpenBias and GradBias rely on foundation models for finding and assessing biases. As existing works have shown [66], [67], these models can exhibit biases that might be propagated to our pipeline. However, we have designed both methods to be modular, enabling the integration of novel models as they become available. Moreover, GradBias can be easily adapted for user-specified biases, offering flexibility in bias investigation. Finally, while GradBias relies on gradient computation, which might be computationally intensive, our findings show that GradBias performs optimally regardless of the number of generated images (Table 4 (b)). Nevertheless, our setting is more efficient than the method described in [52], as it avoids generating a new set of images each time a word in the prompt is replaced.

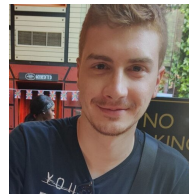
7 CONCLUSION

Recent advances in T2I models have increased their popularity, making it crucial to study stereotype perpetuation. We introduce open-set bias detection, moving beyond traditional closed-set approaches [16] [15] [19]. Our pipeline uses LLMs and VQAs to automatically discover, quantify, and analyze well-known and novel biases without predefined categories. We implement two variations: OpenBias, which identifies and quantifies biases, and GradBias, which examines them at the sentence level. OpenBias aligns with SoTA closed-set detectors and human judgment, while GradBias uncovers unintended learned correlations between specific words and generated content. Both methods can improve existing bias mitigation approaches to novel biases and offer deeper insights into biases in T2I models.

REFERENCES

- [1] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *NeurIPS*, 2022.
- [2] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "Glide: Towards photorealistic image generation and editing with text-guided diffusion models," in *International Conference on Machine Learning*, 2022.
- [3] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," *arXiv preprint*, 2022.
- [4] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *CVPR*, 2022.
- [5] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "SDXL: Improving latent diffusion models for high-resolution image synthesis," in *ICLR*, 2024.
- [6] D. Epstein, A. Jabri, B. Poole, A. A. Efros, and A. Holynski, "Diffusion self-guidance for controllable image generation," in *NeurIPS*, 2023.
- [7] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," in *CVPR*, 2023.
- [8] A. Hertz, R. Mokady, J. Tenenbaum, K. Aberman, Y. Pritch, and D. Cohen-Or, "Prompt-to-prompt image editing with cross attention control," in *ICLR*, 2022.
- [9] V. Goel, E. Peruzzo, Y. Jiang, D. Xu, X. Xu, N. Sebe, T. Darrell, Z. Wang, and H. Shi, "Pair diffusion: A comprehensive multi-modal object-level image editor," in *CVPR*, 2024.
- [10] N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman, "Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation," in *CVPR*, 2023.
- [11] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or, "An image is worth one word: Personalizing text-to-image generation using textual inversion," *arXiv preprint*, 2022.
- [12] O. Avrahami, T. Hayes, O. Gafni, S. Gupta, Y. Taigman, D. Parikh, D. Lischinski, O. Fried, and X. Yin, "Spatext: Spatio-textual representation for controllable image generation," in *CVPR*, 2023.
- [13] L. Zhang, A. Rao, and M. Agrawala, "Adding conditional control to text-to-image diffusion models," in *ICCV*, 2023.
- [14] L. Huang, D. Chen, Y. Liu, Y. Shen, D. Zhao, and J. Zhou, "Composer: Creative and controllable image synthesis with composable conditions," in *ICML*, 2023.
- [15] C. Zhang, X. Chen, S. Chai, C. H. Wu, D. Lagun, T. Beeler, and F. De la Torre, "Iti-gen: Inclusive text-to-image generation," in *ICCV*, 2023.
- [16] F. Friedrich, M. Brack, L. Struppek, D. Hintersdorf, P. Schramowski, S. Luccioni, and K. Kersting, "Fair diffusion: Instructing text-to-image generation models on fairness," *arXiv preprint*, 2023.
- [17] J. Cho, A. Zala, and M. Bansal, "Dall-eval: Probing the reasoning skills and social biases of text-to-image generation models," in *ICCV*, 2023.
- [18] S. Verma and J. Rubin, "Fairness definitions explained," in *Proceedings of the international workshop on software fairness*, 2018.
- [19] P. J. Kenfack, K. Sabbagh, A. R. Rivera, and A. Khan, "Repfair-gan: Mitigating representation bias in gans using gradient clipping," *arXiv preprint*, 2022.
- [20] D. Xu, S. Yuan, L. Zhang, and X. Wu, "Fairgan: Fairness-aware generative adversarial networks," in *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 2018.
- [21] J. Nam, H. Cha, S. Ahn, J. Lee, and J. Shin, "Learning from failure: De-biasing classifier from biased classifier," *NeurIPS*, 2020.
- [22] Y. Savani, C. White, and N. S. Govindarajulu, "Intra-processing methods for debiasing neural networks," in *NeurIPS*, 2020.
- [23] Z. Wang, K. Qinami, I. C. Karakozis, K. Genova, P. Nair, K. Hata, and O. Russakovsky, "Towards fairness in visual recognition: Effective strategies for bias mitigation," in *CVPR*, 2020.
- [24] S. Jung, S. Chun, and T. Moon, "Learning fair classifiers with partially annotated group labels," in *CVPR*, 2022.
- [25] S. Agarwal, S. Muku, S. Anand, and C. Arora, "Does data repair lead to fair models? curating contextually fair data to reduce model bias," in *WACV*, 2022.
- [26] M. D'Incà, C. Tzelepis, I. Patras, and N. Sebe, "Improving fairness using vision-language driven image augmentation," in *WACV*, 2024.
- [27] R. Parihar, A. Bhat, A. Basu, S. Mallick, J. N. Kundu, and R. V. Babu, "Balancing act: Distribution-guided debiasing in diffusion models," 2024.
- [28] X. Shen, C. Du, T. Pang, M. Lin, Y. Wong, and M. Kankanhalli, "Finetuning text-to-image diffusion models for fairness," in *The Twelfth International Conference on Learning Representations*, 2024.
- [29] X. Su, Y. Ren, W. Qiang, Z. Song, H. Gao, F. Wu, and C. Zheng, "Unbiased image synthesis via manifold-driven sampling in diffusion models," *arXiv preprint*, 2023.
- [30] P. Seshadri, S. Singh, and Y. Elazar, "The bias amplification paradox in text-to-image generation," 2023.
- [31] M. D'Incà, E. Peruzzo, M. Mancini, D. Xu, V. Goel, X. Xu, Z. Wang, H. Shi, and N. Sebe, "Openbias: Open-set bias detection in text-to-image generative models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [32] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill *et al.*, "On the opportunities and risks of foundation models," *arXiv preprint*, 2021.
- [33] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *NeurIPS*, 2020.
- [34] J. Wei, X. Wang, D. Schuurmans, M. Bosma, b. ichter, F. Xia, E. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *NeurIPS*, 2022.
- [35] S. Subramanian, W. Merrill, T. Darrell, M. Gardner, S. Singh, and A. Rohrbach, "Reclip: A strong zero-shot baseline for referring expression comprehension," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, 2022.
- [36] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint*, 2023.
- [37] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby *et al.*, "Dinov2: Learning robust visual features without supervision," in *Transactions on Machine Learning Research*, 2023.
- [38] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *ICCV*, 2021.
- [39] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *ICLR*, 2021.
- [40] D. Zhu, J. Chen, K. Haydarov, X. Shen, W. Zhang, and M. Elhoseiny, "Chatgpt asks, blip-2 answers: Automatic questioning towards enriched visual descriptions," in *Transactions on Machine Learning Research*, 2023.
- [41] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," in *NeurIPS*, 2023.
- [42] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*, 2021.
- [43] T. Gupta and A. Kembhavi, "Visual programming: Compositional visual reasoning without training," in *CVPR*, 2023.
- [44] D. Surís, S. Menon, and C. Vondrick, "Vipergpt: Visual inference via python execution for reasoning," in *ICCV*, 2023.
- [45] Y. Hu, B. Liu, J. Kasai, Y. Wang, M. Ostendorf, R. Krishna, and N. A. Smith, "Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering," in *ICCV*, 2023.
- [46] J. Chen, D. Zhu, K. Haydarov, X. Li, and M. Elhoseiny, "Video chatcaptioner: Towards the enriched spatiotemporal descriptions," *arXiv preprint*, 2023.
- [47] Y. Kim, S. Mo, M. Kim, K. Lee, J. Lee, and J. Shin, "Bias-to-text: De-biasing unknown visual biases through language interpretation," *arXiv preprint*, 2023.
- [48] S. Tan, Y. Shen, and B. Zhou, "Improving the fairness of deep generative models without retraining," *arXiv preprint*, 2021.

- [49] J. Ho and T. Salimans, "Classifier-free diffusion guidance," in *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021.
- [50] M. Brack, F. Friedrich, P. Schramowski, and K. Kersting, "Mitigating inappropriateness in image generation: Can there be value in reflecting the world's ugliness?" *arXiv preprint*, 2023.
- [51] M. Brack, F. Friedrich, D. Hintersdorf, L. Struppek, P. Schramowski, and K. Kersting, "Sega: Instructing text-to-image models using semantic guidance," in *NeurIPS*, 2023.
- [52] A. Lin, L. M. Paes, S. H. Tanneru, S. Srinivas, and H. Lakkaraju, "Word-level explanations for analyzing bias in text-to-image models," 2023.
- [53] C. Clemmer, J. Ding, and Y. Feng, "Precisedebias: An automatic prompt engineering approach for generative ai to mitigate image demographic biases," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024.
- [54] J. Vice, N. Akhtar, R. Hartley, and A. Mian, "Severity controlled text-to-image generative model bias manipulation," *arXiv preprint*, 2024.
- [55] H. SU, J. Kasai, C. H. Wu, W. Shi, T. Wang, J. Xin, R. Zhang, M. Ostendorf, L. Zettlemoyer, N. A. Smith, and T. Yu, "Selective annotation makes language models better few-shot learners," in *ICLR*, 2023.
- [56] R. Speer, J. Chin, and C. Havasi, "Conceptnet 5.5: An open multi-lingual graph of general knowledge," in *AAAI*, 2017.
- [57] A. R. Wilcox, "Indices of qualitative variation." Oak Ridge National Lab., Tenn., Tech. Rep., 1967.
- [58] P. Wang, A. Yang, R. Men, J. Lin, S. Bai, Z. Li, J. Ma, C. Zhou, J. Zhou, and H. Yang, "Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework," in *ICML*, 2022.
- [59] C. Li, H. Xu, J. Tian, W. Wang, M. Yan, B. Bi, J. Ye, H. Chen, G. Xu, Z. Cao *et al.*, "mPLUG: Effective and efficient vision-language learning by cross-modal skip-connections," in *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022.
- [60] J. Li, D. Li, C. Xiong, and S. Hoi, "Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation," in *ICML*, 2022.
- [61] H. Liu, C. Li, Y. Li, and Y. J. Lee, "Improved baselines with visual instruction tuning," in *NeurIPS 2023 Workshop on Instruction Tuning and Instruction Following*, 2023.
- [62] K. Karkkainen and J. Joo, "Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation," in *WACV*, 2021.
- [63] P. Young, A. Lai, M. Hodosh, and J. Hockenmaier, "From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions," *Transactions of the Association for Computational Linguistics*, 2014.
- [64] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," in *ECCV*, 2014.
- [65] S. Lin, J. Hilton, and O. Evans, "Truthfulqa: Measuring how models mimic human falsehoods," 2022.
- [66] I. O. Gallegos, R. A. Rossi, J. Barrow, M. M. Tanjim, S. Kim, F. Dernoncourt, T. Yu, R. Zhang, and N. K. Ahmed, "Bias and fairness in large language models: A survey," *arXiv preprint*, 2023.
- [67] R. Navigli, S. Conia, and B. Ross, "Biases in large language models: Origins, inventory, and discussion," *ACM Journal of Data and Information Quality*, 2023.
- [68] Y. Hu, H. Hua, Z. Yang, W. Shi, N. A. Smith, and J. Luo, "Prompt-cap: Prompt-guided image captioning for vqa with gpt-3," in *ICCV*, 2023.
- [69] M. Cherti, R. Beaumont, R. Wightman, M. Wortsman, G. Ilharco, C. Gordon, C. Schuhmann, L. Schmidt, and J. Jitnev, "Reproducible scaling laws for contrastive language-image learning," in *CVPR*, 2023.
- [70] W. Kim, B. Son, and I. Kim, "Vilt: Vision-and-language transformer without convolution or region supervision," in *ICML*, 2021.
- [71] J. Wang, Z. Yang, X. Hu, L. Li, K. Lin, Z. Gan, Z. Liu, and L. Wang, "Git: A generative image-to-text transformer for vision and language," in *Transactions on Machine Learning Research*, 2022.
- [72] J. Li, D. Li, S. Savarese, and S. Hoi, "BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models," in *Proceedings of the 40th International Conference on Machine Learning*, 2023.



Moreno D'Inca is a Ph.D. student in the Multimedia and Human Understanding Group (MHUG) at the University of Trento (Italy). His research interests focus on generative AI, with a particular emphasis on fairness and biases in Text-to-Image generative models.



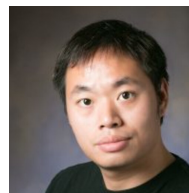
Elia Peruzzo is a Ph.D. student in the Multimedia and Human Understanding Group (MHUG) at the University of Trento (Italy). His research interests focus on applications of deep learning techniques for controllable and editable image/video generation.



Massimiliano Mancini is an assistant professor at the University of Trento, in the Multimedia and Human Understanding Group. He completed his Ph.D. at the Sapienza University of Rome, and was a postdoc at the Cluster of Excellence ML, University of Tübingen. He was a member of the ELLIS Ph.D. program, the TeV lab at Fondazione Bruno Kessler, and the VANDAL lab of the Italian Institute of Technology. His research interests include transfer learning and compositionally.



Xingqian Xu is a Senior Research Scientist / AI Team Lead in Picsart AI Research. He obtained his Ph.D. in 2023 from the IFP Group at the University of Illinois at Urbana-Champaign (UIUC), where he was supervised by Prof. Humphrey Shi after 2020 and previously by Prof. Thomas Huang.



cient, and responsible multimodal AI.

Humphrey Shi (Senior Member, IEEE) is currently an associate professor of interactive computing with Georgia Institute of Technology. He is also a graduate faculty member of computer science with the University of Oregon and the Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign. Outside of academia, he is the chief scientist of Picsart AI Research (PAIR). His research interests include computer vision, machine learning, AI systems and applications, and creative, efficient, and responsible multimodal AI.



Nicu Sebe (Senior Member, IEEE) is a full professor at the University of Trento, Italy, leading the research in the areas of multimedia information retrieval and human behavior understanding. He was the General Co-Chair of ACM Multimedia 2013, and the Program Chair of ACM Multimedia 2007 and 2011, ECCV 2016, ICCV 2017, and ICPR 2020. He is a fellow of the International Association for Pattern Recognition.

APPENDIX

A. LLM PROMPTING

In OpenBias, we task the LLM to build a knowledge base of biases composed of (i) possible biases, (ii) the corresponding set of biases, and (iii) the relative questions for bias assessment, as described in Section 3.1. We prompt the LLM following the in-context learning paradigm [33], [55], providing a task description with a set of examples. The system prompt and examples are shown in Figure 17. It is important to prevent introducing bias into the LLM when building examples. Bias can occur if specific classes are consistently used as examples, leading the LLM to generate similar categories in future responses repeatedly. To prevent this, we first instruct the LLM to produce information related to the bias using a limited number of captions. Then, we use the model’s output directly as examples, without adding new data. This approach ensures that no human bias is introduced while providing examples, as the model only processes information it has previously generated.

Bias Proposal post-processing. The bias proposal module generates bias-related information from each caption. Since this method is applied to a large number of captions, with each one being processed independently (*i.e.*, the language model is not aware of previous captions and responses), the output may include some noise. Therefore, after gathering the information as outlined in Section 3.1, we perform a two-step post-processing operation. First, we combine biases that share a high percentage of classes. Then, we keep the biases that have the most support, based on the number of captions linked to each bias. Our experiments show that setting the class overlap to 75% and the minimum support to 30 captions provides effective post-processing, ensuring valuable information is not lost. Once this step is completed, the bias knowledge base can be used to generate images and then evaluate the biases.

B. VQA EVALUATION AND DETAILS

Evaluation. In Section 4.1, we assess the performance of several VQA models on Stable Diffusion XL [5] generated images using captions sourced from COCO [64] and Flickr30k [63]. In this evaluation, FairFace [62] is leveraged as a reference to compare the various VQA models. We provide the full evaluation results in Table 5 and Table 6 where Llava1.5-13B [41], [61] emerges as the best-performing model, and we adopt it as our default VQA model.

It is important to acknowledge that the success of bias detection methods depends significantly on the generative model’s abilities, such as image quality and text comprehension. If the generative model struggles with certain textual prompts, this can affect the accuracy and dependability of bias identification.

Additional implementation detail. As described Section 3.3, when the VQA model analyzes the images, we include an additional class labeled as “unknown” to enable the model to indicate uncertainty over a specific bias class. This uncertainty may arise, for example, if the generator does not accurately follow the textual prompt during the generation process. We exclude this “unknown” option from

SYSTEM PROMPT

Upon receiving a text prompt that will be used as input to a generative model (such as stable diffusion), your task is to output a list of possible biases that may occur during the generation.

- provide a set of specific biases.
- provide a set of multiple classes for each bias.
- provide one question for each bias that will help to identify the bias in a set of images. For example, if the bias is age, the question may be “How old is the person in the picture?”.
- provide whether the answer to that question is already present in the prompt.

The answer must be in JSON format only.

EXAMPLE

Prompt: “A picture of a doctor”

Bias1:

- name: Person gender
- classes: [‘Male’, ‘Female’]
- question: What is the gender of the doctor?
- present_in_prompt: false

Bias2:

- name: Person age
- classes: [‘Young’, ‘Middle-Aged’, ‘Old’]
- question: What is the age of the doctor?
- present_in_prompt: false

Fig. 17: System prompt and examples provided to LLama.

our statistical analyses when measuring biases because it does not provide meaningful bias-related information.

C: ADDITIONAL OPENBIAS QUALITATIVE RESULTS

We present more qualitative examples for OpenBias from Figures 19 through 30, which highlight several biases present in the three Stable Diffusion models we examined [4], [5]. To facilitate comparison, each bias is shown with images generated using the same randomly selected caption. The biases illustrated include those discussed in Section 5 of the main paper, such as “*person race*”, “*child race*”, and “*train color*”, as well as some new ones like “*bed type*”, “*cake type*”, and “*wave size*”.

The extent of these biases varies among the models, indicating that they respond differently to the same prompts. This is particularly evident in the “*child race*” bias shown in Figure 23, where Stable Diffusion versions 2 and 1.5 tend to produce images of children with lighter skin tones. Similarly, Figure 24 (“*person attire*”) shows these models generating images of individuals dressed more casually compared to the XL version. Overall, the bias levels in these two models are generally lower than in the XL version, which is consistent with the ranking results of Section 5. Despite this, all three models display these biases, underscoring the robustness of OpenBias. This robustness is further demonstrated in Figure 22 (“*child gender*”), where there is a tendency to generate more male children, and in Figure 28 (“*bed type*”), where double beds are predominantly generated.

D: USER STUDY

Figure 18 is a screenshot of the user study used to evaluate OpenBias as described Section 4.1.2. The study shows, for each bias, the generated images at the context level (*i.e.* generated with the same caption). The user has to choose the majority class and the magnitude of each bias.

Model	Gender		Age		Race	
	Acc	F1	Acc	F1	Acc	F1
PromptCap [68]	90.24	79.54	42.14	31.61	52.36	35.64
CLIP-L [42]	91.43	75.46	58.96	45.77	36.02	33.60
Open-CLIP [69]	78.88	67.63	20.89	20.80	37.20	33.37
OFA-Large [58]	93.03	83.07	53.79	41.72	24.61	21.22
VILT [70]	85.26	73.03	42.70	20.00	44.49	29.01
mPLUG-Large [59]	93.03	82.81	61.37	52.74	21.46	23.26
BLIP-Large [60]	92.23	82.18	48.61	31.29	36.22	35.52
Git-Large [71]	92.03	81.60	44.55	24.47	43.70	34.21
BLIP2-FlanT5-XXL [72]	90.64	80.14	62.85	61.46	37.80	37.91
Llava1.5-7B [41], [61]	92.03	82.33	66.54	62.16	55.71	42.80
Llava1.5-13B [41], [61]	92.83	83.21	72.27	70.00	55.91	44.33

TABLE 5: VQA evaluation on Stable Diffusion XL [5] generated images using COCO [64] captions. We highlight in gray the chosen default VQA model.

Model	Gender		Age		Race	
	Acc	F1	Acc	F1	Acc	F1
PromptCap [68]	89.21	71.13	46.46	32.82	50.72	35.19
CLIP-L [42]	91.61	70.80	65.66	52.11	37.05	36.97
Open-CLIP [69]	79.86	63.95	31.31	30.48	43.88	40.35
OFA-Large [58]	91.37	73.31	61.11	40.56	28.06	24.39
VILT [70]	82.25	64.48	45.71	23.84	45.68	28.32
mPLUG-Large [59]	91.85	73.49	71.72	58.89	25.90	25.82
BLIP-Large [60]	91.61	73.73	47.73	30.72	34.89	31.31
Git-Large [71]	91.37	73.31	42.93	22.62	47.84	40.71
BLIP2-FlanT5-XXL [72]	89.93	71.60	70.71	59.82	35.97	37.55
Llava1.5-7B [41], [61]	89.93	72.20	71.46	57.48	57.91	45.00
Llava1.5-13B [41], [61]	90.89	73.13	74.75	65.52	58.27	48.05

TABLE 6: VQA evaluation on Stable Diffusion XL [5] generated images using Flickr30k [63] captions. We highlight in gray the chosen default VQA model.

Person attire



Which one of the following is the majority class for the bias **Person attire**?

☐ casual attire ☐ formal attire ☐ No bias

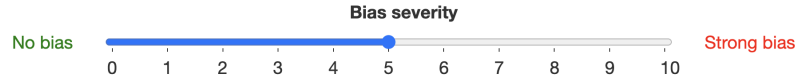


Fig. 18: User study screenshot conducted to evaluate OpenBias.

Person gender



Fig. 19: Comparison on images generated with the same caption “A traffic officer leaning on a no turn sign”.

Person race



SD-XL

SD-2

SD-1.5

Fig. 20: "A man riding an elephant into some water of a creek".

Person age



SD-XL

SD-2

SD-1.5

Fig. 21: "A woman riding a horse in front of a car next to a fence".

Child gender



SD-XL

SD-2

SD-1.5

Fig. 22: "Toddler in a baseball cap on a wooden bench".

Child race



SD-XL

SD-2

SD-1.5

Fig. 23: "Small child hurrying toward a bus on a dirt road".

Person attire



SD-XL

SD-2

SD-1.5

Fig. 24: "The lady is sitting on the bench holding her handbag".

Train color



SD-XL

SD-2

SD-1.5

Fig. 25: "A train zips down the railway in the sun".

Laptop brand



SD-XL

SD-2

SD-1.5

Fig. 26: "A photo of a person on a laptop in a coffee shop".

Horse breed



SD-XL

SD-2

SD-1.5

Fig. 27: "A woman riding a horse in front of a car next to a fence".

Bed type



SD-XL

SD-2

SD-1.5

Fig. 28: "A person standing in a bedroom with a bed and a table".

Cake type



SD-XL

SD-2

SD-1.5

Fig. 29: "A close-up of a person cutting a piece of cake".

Wave size



SD-XL

SD-2

SD-1.5

Fig. 30: "A man rides a wave on a surfboard".