# Seeing is not always believing: Benchmarking Human and Model Perception of AI-Generated Images

**Zeyu Lu**[1,2,*] **Di Huang**[2,3,*] **Lei Bai**[2,*,†] **Jingjing Qu**[2] **Chengyue Wu**[4]
**Xihui Liu**[4] **Wanli Ouyang**[2]

[1] Shanghai Jiao Tong University   [2] Shanghai Artificial Intelligence Laboratory
[3] The University of Sydney   [4] The University of Hong Kong

## Abstract

Photos serve as a way for humans to record what they experience in their daily lives, and they are often regarded as trustworthy sources of information. However, there is a growing concern that the advancement of artificial intelligence (AI) technology may produce fake photos, which can create confusion and diminish trust in photographs. This study aims to comprehensively evaluate agents for distinguishing state-of-the-art AI-generated visual content. Our study benchmarks both human capability and cutting-edge fake image detection AI algorithms, using a newly collected large-scale fake image dataset **Fake2M**. In our human perception evaluation, titled **HPBench**, we discovered that humans struggle significantly to distinguish real photos from AI-generated ones, with a misclassification rate of **38.7%**. Along with this, we conduct the model capability of AI-Generated images detection evaluation **MPBench** and the top-performing model from MPBench achieves a **13%** failure rate under the same setting used in the human evaluation. We hope that our study can raise awareness of the potential risks of AI-generated images and facilitate further research to prevent the spread of false information. More information can refer to https://github.com/Inf-imagine/Sentry.
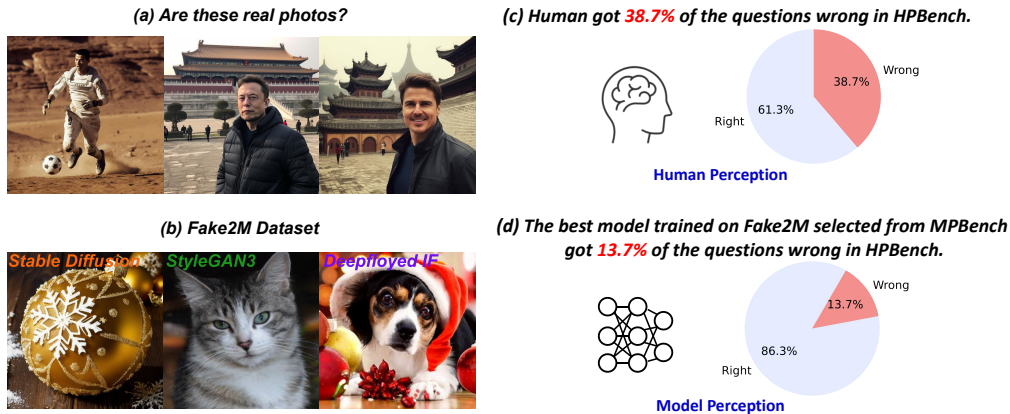
Figure 1: *Left:* (a) Determining whether an image is real is a difficult problem. (b) We introduce a new, large-scale and diverse dataset, named "Fake2M". *Right:* (c) The average human ability to distinguish between high-quality AI-generated images and real images is only 61% in HPBench. (d) The highest performing model trained on the Fake2M dataset, selected from MPBench, achieves an accuracy of 87% on HPBench.

---

*Equal contribution.   †Corresponding author: bailei@pjlab.org.cn.

# 1 Introduction

Photography, which captures images by recording light, has become an integral part of human society, serving as a vital medium for recording real visual information. Its diverse applications range from documenting historical events and scientific discoveries to immortalizing personal memories and artistic expressions. Since the advent of photography with Joseph Nicéphore Niépce's first photograph of photographs in 1826, the global collection of photographs has expanded to an estimated 12.4 trillion photos. This impressive number continues to rise with an annual increase of 1.72 trillion photos (10–14% increasing rate every year) [8].

In contrast to traditional photography, AI-driven image generation harnesses neural networks to learn synthesis rules from extensive image datasets, offering a novel approach for creating high-quality yet fake visual content. By taking image descriptions or random noise as input, AI generates one or several images for users. Early methods employed GANs [14, 31, 39, 41–43, 82] as the generative architecture, while more recent diffusion-based techniques [12, 23, 25, 35, 36, 40, 50, 56, 59, 63–66, 70, 70–73, 83] have showcased improved diversity and generation quality. Therefore, AI-generated contents (AIGC) have gained popularity across various applications, such as AI-assisted design.

Following the rapid advancements in AI-driven image generation algorithms, AI is now capable of producing a wide array of image styles, including those that closely resemble real photographs. At the Sony World Photography Awards in April 2023, the winner Eldagsen refused the award after revealing that his creation is made using AI. Organizers of the award said that Eldagsen had misled them about the extent of AI that would be involved [4]. Consequently, a critical question arises: *Can humans find a reliable solution for distinguishing whether an image is AI-generated?* This inquiry fundamentally questions the reliability of image information for conveying truth. Utilizing generated images to convey false information can have significant social repercussions, often misleading people about nonexistent events or propagating incorrect ideas, as shown in Fig. 1. A recent example involves the circulation of fake images depicting Trump's arrest [3], which garnered substantial social attention.

In order to tackle the challenge of discerning AI-generated images, we carried out an extensive study focusing on human capability and the proficiency of cutting-edge fake image detection AI algorithms. For human capability, we conduct a human evaluation involving participants from diverse backgrounds to determine their ability to distinguish real and fake images. Each participant is tasked with completing a test consisting of 100 randomized questions, where they must decide whether an image is real or generated by AI. As for the evaluation of fake image detection models, existing datasets fell short in supporting this study due to their limited size and inclusion of outdated AI-generated images. Consequently, we assembled a novel dataset, called **Fake2M**, composed of state-of-the-art (SOTA) AI-generated images and real photographs sourced from the internet. Fake2M is a large-scale dataset housing over 2 million AI-generated images from three different models, tailored for training fake image detection algorithms.

For human evaluation of fake image detection task, our key finding underscores that state-of-the-art AI-generated images can indeed significantly deceive human perception. According to our results, participants achieved an average accuracy of only 61.3%, implying a misclassification rate of 38.7%, thus demonstrating the challenge they faced in accurately discerning real images from those generated by AI. Our study also reveals that humans are better at distinguishing AI-generated portrait images compared to other types of AI-generated images.

Turning to the model evaluation for the fake image detection task, our experiments yielded several key insights. Firstly, no single model consistently outperforms the others across all training datasets, suggesting that the optimal model is largely dependent on the specific dataset in use. This underscores the importance of developing a robust model that can consistently achieve superior performance across all training datasets. Secondly, our findings underscored the benefits of diverse training datasets, with models demonstrating improved overall accuracy when exposed to a broader range of visual styles and variations. Lastly, we observed variability in model accuracy when the same model trained on identical datasets was validated using different generation models. This demonstrates that variations in validation dataset generation may influence the performance of model.

In this study, we benchmark the ability of both human and cutting-edge fake image detection AI algorithms. While advancements in image synthesis have enabled individuals to create visually appealing images with quality approaching real photographs, our findings reveal the potential risks associated

with using AIGC for spreading false information and misleading viewers. Our contributions can be listed as follows:

- We introduce a new, large-scale dataset, named "**Fake2M**". To the best of our knowledge, this is the largest and most diverse fake image dataset, designed to stimulate and foster advancements in fake image detection research.
- We establish **HPBench**, a unique benchmark that comprehensively assesses the human capability to discern fake AI-generated images from real ones.
- We establish **MPBench**, an extensive and comprehensive benchmark including 11 fake validation datasets, designed to rigorously assess the model capability for identifying fake images generated by the most advanced generative models currently available.

## 2 Dataset Collection and Generation

### 2.1 Collect Data for Human Evaluation

Table 1: **Number of photographic images used in HPBench across eight categories.**

| Category | Multiperson | Landscape | Man | Woman | Record | Plant | Animal | Object | All |
|---|---|---|---|---|---|---|---|---|---|
| **Number of AI-Generated Images** | 10 | 27 | 17 | 30 | 15 | 13 | 29 | 10 | 151 |
| **Number of Real Images** | 12 | 26 | 44 | 49 | 21 | 18 | 53 | 21 | 244 |

We collect AI-generated images and a set of real photographs across eight categories: Multiperson, Landscape, Man, Woman, Record, Plant, Animal, and Object, as shown Tab. 1.

**Collecting AI-generated photorealistic images.** Firstly, we utilize Midjourney-V5 [7], the state-of-the-art image generation model, to create photorealistic images of the aforementioned eight categories. For each category, we employ diverse prompts to ensure maximum variation. We use specialized prompt suffixes such as "normal color, aesthetic, shocking, very detailed, photographic, 8K, HDR" to improve the authenticity of the images. We notice that in real-world scenarios, people use AI to generate images with the intention of creating high-quality images without visual defects, and users will select the best image from multiple generated images. Therefore, we employ the expertise of annotators to filter out low-quality AI-generated images that can be easily determined as fake photos at the first glance.

**Collecting real photos.** We collect real photos from 500px [1] and Google Images [6] by searching for photos with the same text prompts used for creating AI-generated images in the previous paragraph.

### 2.2 Collect Data for Model Evaluation

Our objective is to investigate whether models can distinguish if an image is AI-generated or not. We constructed 3 training fake datasets with about 2M images, named **Fake2M**, and 11 validation fake datasets with about 257K images using different latest modern generative models, which contain the SOTA Diffusion models (Stable Diffusion [64], IF [5]), the SOTA GAN model (StyleGAN3 [41]), the SOTA autoregressive model (CogView2 [24]), and the SOTA generative model (Midjounrey [7]), as shown in Tab. 2. We describe the details of our datasets in the following subsections.

**Collecting training datasets.** For the text-to-image generation model, we use the first 1M captions from CC3M to generate the corresponding fake images. For the class conditional generation model, we generate the fake images directly. We describe the specific settings of the 3 training fake datasets in Tab 2, as follows: **(1) "SD-V1.5Real-dpms-25":** Stable Diffusion v1.5 Realistic Vision V2.0 [10] (the top popular models in CIVITAI) with DPM-Solver 25 steps to generate the corresponding fake images. **(2) "IF-V1.0-dpms++-25":** IF v1.0 [5] with DPM-Solver++ [49] 25 steps to generate the corresponding fake images. **(3) "StyleGAN3":** To match the training datasets domain in StyleGAN3 [41] and to maximize the diversity of the model, we use StyleGAN3-t-ffhq to generate the 35K fake images, StyleGAN3-r-ffhq to generate the 35K fake images, StyleGAN3-t-metfaces to generate the 650 fake images, StyleGAN3-r-metfaces to generate the 650 fake images, stylegan3-t-afhqv2 to generate the 8K fake images and stylegan3-r-afhqv2 to generate the 8K fake images.

Table 2: **Detailed information of the datasets used in MPBench. R** denotes the dataset consisting entirely of real images. **F** denotes the dataset consisting entirely of fake images. ✔ denotes existing datasets. ✗ denotes the datasets provided in this work. "Diff" refers to diffusion model, "AR" refers to autoregressive model and "Unk." refers to unknown model.

| Dataset | CC3M-Train | StyleGAN3-Train | SD-V1.5Real-dpms-25 | IF-V1.0-dpms++-25 | StyleGAN3 | ImageNet-Test | CelebA-HQ-Train | CC3M-Val | SD-V2.1-dpms-25 | SD-V1.5-dpms-25 | SD-V1.5Real-dpms-25 | IF-V1.0-dpms++-10 | IF-V1.0-dpms++-25 | IF-V1.0-dpms++-50 | IF-V1.0-ddim-50 | IF-V1.0-ddpms-50 | Cogview2 | Midjourney | StyleGan3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Category | Train | | | | | Validate | | | | | | | | | | | | | |
| | R | R | F | F | F | R | R | R | F | F | F | F | F | F | F | F | F | F | F |
| Generator | - | - | Diff. | Diff. | GAN | - | - | - | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | AR | Unk. | GAN |
| Numbers | 1M | 87K | 1M | 1M | 87K | 100K | 24K | 15K | 15K | 15K | 15K | 15K | 15K | 15K | 15K | 15K | 22K | 5.5K | 60K |
| This work | ✗ | ✗ | ✔ | ✔ | ✔ | ✗ | ✗ | ✗ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |

In order to align with the domain of former fake datasets, we used the corresponding real images datasets and the specific settings of the 2 training real datasets in Tab 2, as follows: **(1) "CC3M-Train":** To conform with the former fake datasets, we use the corresponding real images from CC3M [68]. **(2) "StyleGAN3-Train":** To match the former fake dataset, we use the StyleGAN3 training datasets: FFHQ [42], AFHQv2 [19] and MetFaces [41].

**Collecting validation datasets.** For the text-to-image generation model, we use the whole 15K captions from CC3M validation dataset to generate the corresponding fake images. For the class conditional generation model, we generate fake images directly. For Midjourney, we crawled 5.5K images from the community as a validation set. We describe the specific settings of the 11 validation fake datasets in Tab 2, as follows: **(1) "SD-V2.1-dpms-25":** Stable Diffusion v2.1 with DPM-Solver [48] 25 steps to generate the corresponding fake images. **(2) "SD-V1.5-dpms-25":** Stable Diffusion v1.5 with DPM-Solver 25 steps to generate the corresponding fake images. **(3) "SD-V1.5Real-dpms-25":** Stable Diffusion v1.5 Realistic Vision V2.0 with DPM-Solver 25 steps to generate the corresponding fake images. **(4) "IF-V1.0-dpms++-10":** IF v1.0 with DPM-Solver++ 10 steps to generate the corresponding fake images. **(5) "IF-V1.0-dpms++-25":** IF v1.0 with DPM-Solver++ 25 steps to generate the corresponding fake images. **(6) "IF-V1.0-dpms++-50":** IF v1.0 with DPM-Solver++ 50 steps to generate the corresponding fake images. **(7) "IF-V1.0-ddim-50":** IF v1.0 with DDIM [70] 50 steps to generate the corresponding fake images. **(8) "IF-V1.0-ddpm-50":** IF v1.0 with DDPM [36] 50 steps to generate the corresponding fake images. **(9) "Cogview2":** Cogview2 to generate the corresponding fake images. **(10) "Midjourney":** We crawl 5K images from the community as a validation set. **(11) "StyleGan3":** To maximize the diversity of the model, we generate 10K images each using StyleGAN3-t-ffhq, StyleGAN3-r-ffhq, StyleGAN3-t-metfaces, StyleGAN3-r-metfaces, stylegan3-t-afhqv2 and stylegan3-r-afhqv2, for a total of 60K images.

We use 3 common real datasets as our validation dataset and the specific settings as follows: **(1) "ImageNet-Test":** The test dataset of ImageNet1k [22]. **(2) "CelebA-HQ-Train":** The training dataset of CelebA-HQ [44]. **(3) "CC3M-Val":** We use the validation dataset of CC3M as our validation dataset.

## 2.3 Comparison with Other Datasets

As shown in Tab. 3, our dataset is the largest and most diverse public general fake image dataset, designed to stimulate and foster advancements in fake image detection research.

Table 3: **Comparison with other fake image datasets.** "Content Type" means the type of the content in each dataset ("Face" means that the content in this dataset is mostly faces, such as FFHQ [42]. "Object" means that the content in this dataset is mainly composed of a limited number of objects, such as ImageNet [22]. "General" means that the content in this dataset is general, not limited to some objects, faces or art, such as CC3M [68]). "Generator Type" means the type of generator used in our dataset. "Public" means the dataset is publicly accessible. "Fake Image Number" represents the number of fake images provided by this dataset.

| Dataset | Content Type | Generator Type | | | Public | Generator Number | Fake Image Number |
| | | GAN | Diffusion | AutoRegressive | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| FakeSpotter [78] | Face | ✔ | ✗ | ✗ | ✗ | 7 | 5K |
| DFFD [21] | Face | ✔ | ✗ | ✗ | ✔ | 8 | 240K |
| APFDD [29] | Face | ✔ | ✗ | ✗ | ✗ | 1 | 5K |
| ForgeryNet [34] | Face | ✔ | ✗ | ✗ | ✔ | 15 | 1.4M |
| DeepArt [80] | Art | ✗ | ✔ | ✗ | ✔ | 5 | 73K |
| CNNSpot [79] | Object | ✔ | ✗ | ✗ | ✔ | 11 | 362K |
| IEEE VIP Cup [76] | Object | ✔ | ✔ | ✗ | ✗ | 5 | 7K |
| CiFAKE [13] | Object | ✗ | ✔ | ✗ | ✔ | 1 | 60K |
| **Ours** | **General** | ✔ | ✔ | ✔ | ✔ | **11** | **2.3M** |

# 3   HPBench: Human Perception of AI-Generated Images Evaluation

Our objective is to investigate whether humans can distinguish if an image is AI-generated or not. We collect a set of AI-generated photorealistic images and real photographs, and conduct a human evaluation on the collected images to build **HPBench**. We describe the details of our data collection, human evaluation, and metrics for analyzing the results in the following subsections.

## 3.1   Evaluation Setup

Table 4: **Number of participants across different backgrounds.** "w/ AIGC" refers to participants who have played with AIGC. "w/o AIGC" refers to participants who have not played with AIGC.

| | Gender | | Background | | Age | | All |
| **Category** | Male | Female | w/ AIGC | w/o AIGC | 20~29 | 30~45 | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| **Number** | 31 | 19 | 27 | 23 | 42 | 8 | 50 |

**A high-quality fifty-participant human evaluation.** In order to ensure comprehensiveness, fairness, and quality of the evaluation, we recruit a total of 50 participants to participate in our human evaluation instead of using crowdsourcing and make efforts to ensure the diversity of the participants, as shown in Tab. 4. Each participant is asked to complete a questionnaire consisting of 100 questions to determine whether the image is generated by AI or not without any time limit.

Table 5: **Human evaluation metrics across nine categories using all participant data.**

| | Category | | | | | | | | |
| **Metric** | All | Multiperson | Landscape | Man | Woman | Record | Plant | Animal | Object |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| **Accuracy** ↑ | 0.6134 | 0.6750 | 0.5650 | 0.6433 | 0.6637 | 0.6233 | 0.5983 | 0.6133 | 0.5083 |
| **Precision** ↑ | 0.6278 | 0.7075 | 0.5657 | 0.6666 | 0.6765 | 0.6340 | 0.6213 | 0.6156 | 0.5112 |
| **Recall** ↑ | 0.5577 | 0.5966 | 0.5714 | 0.5733 | 0.6275 | 0.5833 | 0.5033 | 0.6033 | 0.3800 |
| **FOR** ↓ | 0.3981 | 0.3487 | 0.4358 | 0.3742 | 0.3473 | 0.3858 | 0.4173 | 0.3888 | 0.4933 |

**Evaluation metrics.** We employ four commonly used evaluation metrics to analyze our results and highlight their respective meanings in the context of our problem. We define positive samples as AI-generated images and negative samples as real images for our problem, and then calculate Accuracy, Precision, Recall, and False Omission Rate (FOR) in the context of our problem in Tab. 5.

### 3.2 Results and Analysis

#### 3.2.1 Overall Ability to Distinguish Real and AI-generated Images

**Results.** The results are shown in Fig. 2. Our study indicates that participants on average are able to correctly distinguish 61.3% of the images, while 38.7% of the images are misclassified. The highest-performing participant is able to correctly distinguish 73% of the images, while the lowest-performing participant is only able to correctly distinguish 40% of the images. These observations demonstrate that a combination of real and AI-generated fake images can easily deceive people. Moreover, the results reveal that humans have an average probability of 66.9% of correctly identifying real photos from the internet, whereas for AI-generated images, people are more likely to be misled and incorrectly identify them as real with an average probability of 44.2%.

**Analysis.** An intriguing observation from the above data is that even for real images collected from the Internet, participants are only able to correctly identify 66.9% of the images that should have been correctly sorted out 100% of the time. This finding demonstrates that AI-generated images not only convey incorrect information to humans but also erode people's trust in accurate
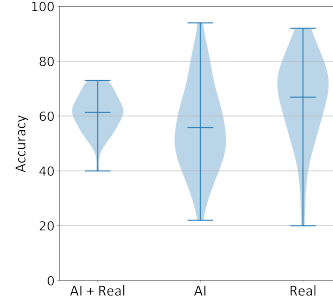


Figure 2: **Human evaluation score distribution using all participant data.** "AI" represents only counting the AI-generated images. "Real" represents only counting the real images. "AI + Real" represents counting all the images.

information. Additionally, our study reveals that humans possess a superior ability to identify real photos than fake photos, which can be attributed to their lifetime experience. Despite this, the difference between real image perception and AI-generated image perception is only 11.1%, and we believe that future generative AI models could further narrow this gap.

#### 3.2.2 Distinguishing Abilities of Participants with Various Personal Backgrounds

**Results.** We explore the effect of AIGC background (refers to participants who have played with AIGC) in our human evaluation. In Fig. 3, we can see that individuals with AIGC background scored slightly better than those without AIGC background (+2.7%). Interestingly, their AIGC background seem to have slight effects on their ability to identify real images (+0.7%). However, when it comes to AI-generated images, participants with AIGC background performs significantly better, with a boost of 3.7%.

**Analysis.** Our research and analysis demonstrate that knowledge and experience with AIGC do not play a significant role in their ability to distinguish between real and AI-generated images. Specifically, when it comes to AI-generated images, participants with AIGC background have a slightly improved ability to distinguish between the two categories. Furthermore, this suggests that additional training and exposure to generative models may be beneficial for individuals by helping them make more informed decisions and avoid any potential risks from fake images.
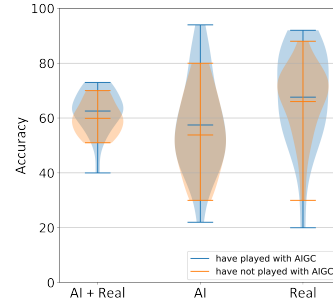


Figure 3: **Human evaluation score distribution calculated by data from participants with AIGC background and without AIGC background.**

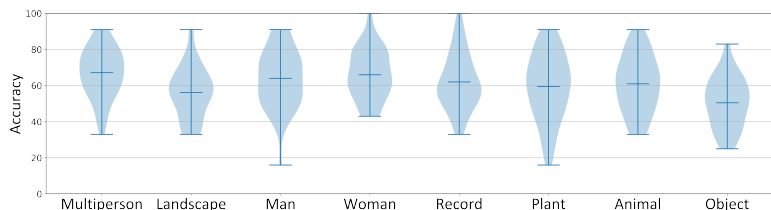#### 3.2.3 Distinguishability of different photo categories



Figure 4: **Human evaluation score distribution across eight categories using all participant data.**

**Results.** We aim to investigate the difficulties for humans to distinguish real and fake images of different image categories. The results are presented in Fig. 4 and Tab. 5. It is observed that the participants had varying ability levels in distinguishing real or fake images from different categories. For instance, the participants can distinguish real or fake images from the category *Multiperson* with a relatively high accuracy rate of 67.5%, whereas they can only correctly distinguish real or fake from the category *Object* with a much lower accuracy rate of 50.8%. Our results indicate a significant 16.7% difference between the accuracy rates of the easiest and most challenging categories. These findings suggest that there may be significant differences in the way humans perceive different photo categories, which could have important implications for further research in this area. It also indicates that the current AI-based generative models may be good at generating some categories but not so good at generating other categories.

**Analysis.** The presented data offers valuable insights into the strengths and limitations of generative AI. The study indicates that current state-of-the-art AI image generation models excel in creating *Object* images that are incredibly realistic. However, AI still struggles when it comes to generating human images. In Tab. 5, *Multiperson*, *Man*, and *Woman* categories are easier to be distinguished as real or fake images compared to others. This phenomenon may be attributed to the fact that the human is more sensitive to images of humans, which is an essential aspect of our cognitive processing.

Table 6: **Statics of judgement criteria from participants correctly identifying fake images.**

| Category | Detail | Smooth | Blur | Color | Shadow & Light | Daub | Rationality | Intuition |
|---|---|---|---|---|---|---|---|---|
| **Number** | 332 | 205 | 142 | 122 | 95 | 59 | 57 | 169 |
| **Percent** | 28% | 17% | 12% | 10% | 8% | 5% | 5% | 14% |

#### 3.2.4 Results of the Judgment Criteria and Analysis of the AIGC Defects

**Results.** We predefined eight judgment criteria options based on our experience and the specific statistics of judgment criteria for correctly selecting the AI-generated images are shown in Tab. 6. The most common issues are "Detail problem" and "Smooth problem" with a high rate of 28% and 17%, respectively. It is also worth noting that the proportion of participants who choose "Intuition" is about 14%, indicating that it is difficult for people to describe obvious defects in AI-generated images, even if they can successfully identify the AI-generated images.

**Analysis.** Our experiment has revealed that even high-quality AI-generated images still exhibit various imperfections and shortcomings. Furthermore, the current SOTA image generation model has limited capability in generating fine details and often generates portraits that are overly smoothed.

## 4 MPBench: Model Perception of AI-Generated Images Evaluation

Our objective is to investigate whether AI models can distinguish if an image is AI-generated or not. We conduct a large and comprehensive model evaluation using **Fake2M** training dataset and 11 fake validation datasets to build **MPBench**.

### 4.1 Experiments Setup

We compare the following SOTA methods: (1) Wang *et al.* [79] proposed to finetune a classification network to give a real/fake decision for an image using *Blur* and *JPEG* augmentation with a pre-trained visual backbone, such as ResNet-50 [33] and ConvNext-S [47]. (2) Ojha *et al.* [57] proposed to use a frozen CLIP-ViT-L [61] for backbone and train the last linear layer for the same task [57].

To ensure a fair comparison, we train each model using four different dataset settings, consisting of 1 million fake images and an equivalent number of real images: (1) Dataset Setting A: "SD-V1.5Real-dpms-25" dataset with 1M fake images and "CC3M-Train" dataset with 1M real images. (2) Dataset Setting B: "IF-V1.0-dpms++-25" dataset with 1M fake images and "CC3M-Train" dataset with 1M real images. (3) Dataset Setting C: "StyleGAN3" dataset with 87K fake images and "StyleGAN3-Train" with 87K real images. (4) Dataset Setting D: "SD-V1.5Real-dpms-25, IF-V1.0-dpms++-25, StyleGAN3, CC3M-Train, StyleGAN3-Train" dataset with 460K fake images selected from "SD-V1.5Real-dpms-25" dataset, 460K fake images selected from "IF-V1.0-dpms++-25" dataset, 87K

Table 7: **Quantitative comparison of five models under four training dataset settings with fourteen validation datasets.** "Diff" refers to diffusion model, "AR" refers to autoregressive model and "Unk." refers to unknown model. **Real (R)** denotes the dataset consisting entirely of real images. **Fake (F)** denotes the dataset consisting entirely of fake images.

| Model | Training Dataset | ImageNet-Test | CelebA-HQ-Train | CC3M-Val | Average Acc. | SD-V2.1-dpm-25 | SD-V1.5-dpm-25 | SD-V1.5Real-dpm-25 | IF-V1.0-dpm++-10 | IF-V1.0-dpm++-25 | IF-V1.0-dpm++-50 | IF-V1.0-ddim-50 | IF-V1.0-ddpm-50 | Cogview2 | Midjourney | StyleGAN3 | Average Acc. | Total Average Acc. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | **Real** | | | | **Fake** | | | | | | | | | | | | **R+F** |
| | | - | - | - | - | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | AR | Unk. | GAN | - | - |
| ConvNext-S(B+J 0.1) | | 95.3 | 99.9 | 99.9 | 98.3 | 48.6 | 99.9 | 85.1 | 92.9 | 57.6 | 55.9 | 41.1 | 72.6 | 55.7 | 41.6 | 37.4 | 62.5 | 70.2 |
| ConvNext-S(B+J 0.5) | Dataset Setting A: | 95.9 | 99.9 | 99.9 | _98.5_ | 53.5 | 100 | 83.3 | 91.1 | 50.2 | 49.6 | 35.3 | 66.9 | 54.9 | 44.7 | 35.7 | 60.4 | 68.6 |
| ResNet50(B+J 0.1) | SD-V1.5Real-dpm-25 (1M) | 93.0 | 95.3 | 95.6 | 94.6 | 71.7 | 71.8 | 98.6 | 57.2 | 26.6 | 29.0 | 23.6 | 47.9 | 11.8 | 40.3 | 9.1 | 44.3 | 55.1 |
| ResNet50(B+J 0.5) | CC3M-Train (1M) | 93.2 | 95.5 | 95.8 | 94.8 | 70.6 | 70.4 | 98.5 | 52.1 | 23.6 | 26.1 | 21.6 | 46.2 | 10.6 | 36.6 | 7.2 | 42.1 | 53.4 |
| CLIP-ViT-L(LC) | | 49.6 | 87.0 | 75.4 | 70.6 | 73.3 | 86.6 | 97.6 | 93.9 | 77.8 | 71.4 | 84.1 | 90.9 | 86.6 | 88.7 | 86.1 | **85.1** | 82.0 |
| ConvNext-S(B+J 0.1) | | 87.6 | 99.9 | 99.9 | _95.8_ | 2.2 | 34.5 | 2.5 | 99.7 | 99.9 | 99.9 | 11.1 | 66.4 | 19.2 | 8.9 | 10.1 | 41.3 | 52.9 |
| ConvNext-S(B+J 0.5) | Dataset Setting B: | 87.8 | 99.9 | 99.9 | _95.8_ | 3.9 | 39.1 | 3.9 | 99.6 | 99.9 | 99.8 | 18.5 | 79.2 | 25.8 | 8.1 | 8.0 | 44.1 | 55.2 |
| ResNet50(B+J 0.1) | IF-V1.0-dpms++-25 (1M) | 89.4 | 95.8 | 95.0 | 93.4 | 37.5 | 56.5 | 20.0 | 84.0 | 95.6 | 91.7 | 39.7 | 69.4 | 45.3 | 15.7 | 8.8 | 51.2 | 60.3 |
| ResNet50(B+J 0.5) | CC3M-Train (1M) | 90.8 | 95.6 | 94.5 | 93.6 | 41.6 | 58.8 | 21.7 | 82.3 | 95.2 | 91.3 | 47.0 | 79.7 | 56.1 | 18.3 | 6.8 | _54.4_ | 62.8 |
| CLIP-ViT-L(LC) | | 81.5 | 83.3 | 93.0 | 85.9 | 38.2 | 18.5 | 13.1 | 80.4 | 79.7 | 70.9 | 61.1 | 77.7 | 76.6 | 33.7 | 32.8 | 52.9 | 60.0 |
| ConvNext-S(B+J 0.1) | | 58.7 | 81.0 | 62.7 | 67.4 | 44.9 | 42.7 | 39.4 | 40.7 | 42.3 | 42.2 | 35.7 | 39.9 | 53.2 | 41.6 | 69.9 | 44.7 | 49.6 |
| ConvNext-S(B+J 0.5) | Dataset Setting C: | 67.7 | 92.6 | 71.2 | _77.1_ | 33.1 | 33.7 | 32.1 | 33.6 | 36.5 | 34.2 | 35.0 | 35.5 | 44.8 | 28.5 | 42.6 | 35.4 | 44.3 |
| ResNet50(B+J 0.1) | StyleGAN3 (87K) | 31.5 | 14.3 | 39.1 | 28.3 | 67.7 | 62.7 | 68.1 | 65.3 | 71.0 | 67.0 | 58.8 | 60.9 | 60.9 | 81.9 | 90.9 | 68.6 | 60.0 |
| ResNet50(B+J 0.5) | StyleGAN3-Train (87K) | 70.6 | 29.5 | 63.4 | 54.5 | 37.3 | 34.9 | 40.0 | 43.3 | 43.9 | 41.9 | 34.8 | 38.4 | 47.8 | 58.1 | 81.3 | 45.6 | 47.5 |
| CLIP-ViT-L(LC) | | 37.6 | 85.9 | 70.2 | 64.5 | 85.7 | 92.9 | 94.8 | 95.3 | 89.5 | 84.0 | 87.5 | 93.0 | 84.6 | 81.4 | 61.9 | **86.4** | 81.7 |
| ConvNext-S(B+J 0.1) | Dataset Setting D: | 95.4 | 99.9 | 99.9 | 98.4 | 39.9 | 95.1 | 99.9 | 99.7 | 99.8 | 99.7 | 43.1 | 90.4 | 48.3 | 32.4 | 99.8 | 77.1 | 81.6 |
| ConvNext-S(B+J 0.5) | SD-V1.5Real-dpms-25 (460K) | 97.7 | 99.9 | 99.9 | **99.1** | 52.8 | 92.9 | 99.9 | 99.5 | 99.7 | 99.3 | 46.4 | 91.7 | 48.6 | 35.0 | 99.9 | 78.7 | **83.0** |
| ResNet50(B+J 0.1) | IF-V1.0-dpms++-25 (460K) | 77.3 | 93.6 | 91.7 | 87.5 | 83.9 | 90.4 | 96.8 | 91.2 | 90.6 | 86.7 | 52.3 | 82.2 | 56.4 | 49.4 | 80.8 | 78.2 | 80.2 |
| ResNet50(B+J 0.5) | StyleGAN3 (87K) CC3M-Train (1M) | 84.5 | 94.7 | 92.1 | 90.4 | 84.4 | 89.4 | 96.1 | 88.9 | 89.4 | 85.6 | 52.4 | 83.8 | 55.4 | 45.8 | 69.5 | 76.4 | 79.4 |
| CLIP-ViT-L(LC) | StyleGAN3-Train (87K) | 50.4 | 96.0 | 79.8 | 75.4 | 66.9 | 87.8 | 93.6 | 96.1 | 85.9 | 78.9 | 84.9 | 92.6 | 87.4 | 79.8 | 62.8 | _83.3_ | 81.6 |

fake images from "StyleGAN3" dataset, 1M real images from "CC3M-Train" dataset and 87K real images from "StyleGAN3-Train" dataset.

We eval all the models using 11 validation datasets to build **MPBench**, as shown in Tab 7. To clearly show the performance of the models under different validation datasets, we directly report the classification accuracy of each model.

## 4.2 Results and Analysis

### 4.2.1 Comparative Analysis of Accuracy Across Various Models

**Results.** We conducted a systematic analysis of the mean accuracy achieved by different models, each trained using an identical dataset. Our findings are summarized in Tab. 7. Notably, the leading performing models for fake image detection vary depending on the training dataset. For example, under Dataset Setting A, CLIP-ViT-L(LC) excels, whereas under Dataset Setting B, ResNet50 (B+J 0.5) outperforms others. Furthermore, CLIP-ViT-L(LC) achieves the highest fake image detection accuracies of 86.4% and 83.3% under Dataset Setting C and Dataset Setting D respectively. A compelling finding from our study is the absence of a single model that consistently delivers superior performance across all dataset settings for fake image detection. However, our experiments also suggest that ConvNext models are more adept at achieving higher average accuracies in detecting real images. In all four dataset settings, ConvNext consistently achieves top accuracies—98.5%, 95.8%, 77.1%, and 99.1% for Dataset Setting A,B,C,D, respectively.

**Analysis.** Our observations highlight the need for models that perform optimally on both real and fake images. According to our study, while CLIP-ViT-L(LC) most often produces the best results for fake images, ConvNext outperforms CLIP-ViT-L(LC) in real image detection. However, real-world applications necessitate discerning between real and fake images. This underscores the need for future research to strike a balance between detection capabilities for both categories.

### 4.2.2 Comparative Analysis of Accuracy Across Various Training Datasets

**Results.** We conduct an empirical analysis of the average accuracy of various models, each trained using different dataset configurations. The results, as presented in Tab. 7, indicate that Dataset Setting D generally yields superior model performance. ConvNext and ResNet50 models specifically show marked improvements when trained on Dataset Setting D. Meanwhile, CLIP-ViT-L (LC) demonstrates comparable performance across different dataset configurations, attaining accuracies of 82.0% with Dataset Setting A and 81.6% with Dataset Setting D.

**Analysis.** Our experimental findings suggest that the choice of training dataset significantly influences model performance in fake image detection tasks. Notably, a diversified dataset like Dataset Setting D, which includes five distinct generative models, seems to enhance overall model accuracy. This improvement is likely due to the exposure to a broader spectrum of generative styles and variations that a diversified dataset offers. Additionally, our results highlight the aptitude of the proposed dataset, **Fake2M**, which features diverse data sources, for more generalized fake image detection.

### 4.2.3 Comparative Analysis of Accuracy Across Various Validation Datasets

**Results.** In our experiment, we assessed the accuracy of a trained model across various validation datasets. As depicted in Tab. 7, the performance of a model varies significantly based on the generation models, sampling methods, and sampling steps used in the validation set. For instance, the ConvNext-S (B+J 0.5) model, when trained using Dataset Setting D, displayed a broad range of validation results for fake image detection, with accuracies spanning between 35.0% and 99.9%. These findings underscore the influence of validation set characteristics on a model's performance in the task of fake image detection.

**Analysis.** In realistic applications, fake images can originate from a variety of generation models with a diverse range of hyperparameters. Consequently, an ideal fake image detection model should demonstrate consistent proficiency across this broad spectrum of generation settings. However, our experimental results show that existing models do not meet this requirement, indicating a need for the development of new approaches to handle the challenge of detecting fake images generated under varying settings.

### 4.2.4 Evaluate the best model under the same setting used in HPBench.

Based on its highest total average accuracy in MPBench, we select ConvNext-S (B+J 0.5) with Dataset Setting D as the best model. We then evaluate this model under the same setting in HPBench which consists of 50 real images and 50 fake images for testing and it achieves a 13% failure rate.

## 5 Related Work

**Image Generation Models.** State-of-the-art text-to-image synthesis approaches such as DALL·E 2 [63], Imagen [66], Stable Diffusion [64], IF [5], and Midjourney [7] have demonstrated the possibility of that generating high-quality, photorealistic images with diffusion-based generative models trained on large datasets [67]. Those models have surpassed previous GAN-based models [14, 31, 42, 82] in both fidelity and diversity of generated images, without the instability and mode collapse issues that GANs are prone to. In addition to diffusion models, other autoregressive models such as Make-A-Scene [27], CogView [24], and Parti [81] have also achieved amazing performance.

**Fake Images Generation and Detection.** In recent years, there have been many works [11, 16, 20, 26, 52, 55, 57, 79, 84] exploring how to distinguish whether an image is AI-generated. These works focus on fake contents generated by GANs or small generation models [14, 31, 42, 82]. Due to the limited quality of images generated by those methods, it is easy for humans to distinguish whether a photo is AI-generated or not. However, as the quality of generated images continues to improve with the advancement of recent generative models [7, 63, 64, 66], it has become increasingly difficult for humans to identify whether an image is generated by AI or not. The need for detecting fake images has existed even before we had powerful image generators.

# 6 Conclusion

In this study, we present a comprehensive evaluation of both human discernment and contemporary AI algorithms in detecting fake images. Our findings reveal that humans can be significantly deceived by current cutting-edge image generation models. In contrast, AI fake image detection algorithms demonstrate a superior ability to distinguish authentic images from fakes. Despite this, our research highlights that existing AI algorithms, with a considerable misclassification rate of 13%, still face significant challenges. We anticipate that our proposed dataset, **Fake2M**, and our dual benchmarks, **HPBench** and **MPBench**, will invigorate further research in this area and assist researchers in crafting secure and reliable AI-generated content systems. As we advance in this technological era, it is crucial to prioritize responsible creation and application of generative AI to ensure its benefits are harnessed positively for society.

# References

[1] 500px. https://500px.com/. Accessed: 2023-04-17. 3

[2] bbc news: "art is dead dude" - the rise of the ai artists stirs debate. https://www.bbc.com/news/technology-62788725. Accessed: 2023-04-18. 29

[3] bbc news: Fake trump arrest photos: How to spot an ai-generated image. https://www.bbc.com/news/world-us-canada-65069316. Accessed: 2023-04-18. 2, 29

[4] bbc news: Sony world photography award 2023: Winner refuses award after revealing ai creation. https://www.bbc.com/news/entertainment-arts-65296763. Accessed: 2023-04-17. 2

[5] Deepfloyd. if. https://github.com/deep-floyd/IF. Accessed: 2023-06-7. 3, 9, 17, 30

[6] Google images. https://images.google.com/. Accessed: 2023-04-17. 3

[7] Midjourney. https://www.midjourney.com/. Accessed: 2023-04-17. 3, 9, 25, 26

[8] photutorial: Number of photos (2023): Statistics, facts, & predictions. https://photutorial.com/. Accessed: 2023-04-18. 2

[9] Rembrandt's the night watch painting restored by ai. https://www.bbc.com/news/technology-57588270. Accessed: 2023-04-18. 29

[10] Stable diffusion v1.5. realistic vision v2.0. https://civitai.com/models/4201/realistic-vision-v20. Accessed: 2023-06-7. 3

[11] Vishal Asnani, Xi Yin, Tal Hassner, Sijia Liu, and Xiaoming Liu. Proactive image manipulation detection. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 9, 26

[12] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 2

[13] Jordan J Bird and Ahmad Lotfi. Cifake: Image classification and explainable identification of ai-generated synthetic images. *arXiv preprint arXiv:2303.14126*, 2023. 5

[14] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *International Conference on Learning Representations*, 2019. 2, 9, 25, 26

[15] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 2020. 25

[16] Lucy Chai, David Bau, Ser-Nam Lim, and Phillip Isola. What makes fake images detectable? understanding properties that generalize. *European Conference on Computer Vision*, 2020. 9, 24, 26

[17] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for contrastive language-image learning. *arXiv preprint arXiv:2212.07143*, 2022. 28

[18] Bobby Chesney and Danielle Citron. Deep fakes: A looming challenge for privacy, democracy, and national security. *California Law Review*, 2019. 26

[19] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 4, 16

[20] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva. Splicebuster: A new blind image splicing detector. *IEEE International Workshop on Information Forensics and Security*, 2015. 9, 26

[21] Hao Dang, Feng Liu, Joel Stehouwer, Xiaoming Liu, and Anil K. Jain. On the detection of digital face manipulation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 5

[22] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. 2009. 4, 5, 18

[23] Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*. 2

[24] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. Cogview: Mastering text-to-image generation via transformers. *Advances in Neural Information Processing Systems*, 2021. 3, 9, 25

[25] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. *arXiv preprint arXiv:2302.03011*, 2023. 2

[26] Joel Frank, Thorsten Eisenhofer, Lea Schonherr, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. Leveraging frequency analysis for deep fake image recognition. *International Conference on Machine Learning*, 2020. 9, 26

[27] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. *European Conference on Computer Vision*, 2022. 9, 25

[28] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. 29

[29] Apurva Gandhi and Shomik Jain. Adversarial perturbations fool deepfake detectors. *International Joint Conference on Neural Networks*, 2020. 5

[30] Robert Geirhos, Kantharaju Narayanappa, Benjamin Mitzkus, Tizian Thieringer, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Partial success in closing the gap between human and machine vision. *Advances in Neural Information Processing Systems*, 2021. 21

[31] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014. 2, 9, 25, 26

[32] Teddy Surya Gunawan, Siti Amalina Mohammad Hanafiah, Mira Kartiwi, Nanang Ismail, Nor Farahidah Za'bah, and Anis Nurashikin Nordin. Development of photo forensics algorithm by detecting photoshop manipulation using error level analysis. *Indonesian Journal of Electrical Engineering and Computer Science*, 2017. 24

[33] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016. 7, 28

[34] Yinan He, Bei Gan, Siyu Chen, Yichun Zhou, Guojun Yin, Luchuan Song, Lu Sheng, Jing Shao, and Ziwei Liu. Forgerynet: A versatile benchmark for comprehensive forgery analysis. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 5

[35] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey A. Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2

[36] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 2020. 2, 4, 16, 25

[37] Justin M. Johnson and Taghi M. Khoshgoftaar. Survey on deep learning with class imbalance. *Journal of Big Data*, 2019. 30

[38] Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung Park. Scaling up gans for text-to-image synthesis. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 30

[39] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *International Conference on Learning Representations*, 2018. 2

[40] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 2022. 2

[41] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 2021. 2, 3, 4, 16, 17, 30

[42] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019. 4, 5, 9, 16, 18, 25, 26

[43] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 2

[44] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 4

[45] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *IEEE/CVF International Conference on Computer Vision*, 2021. 26

[46] Zhengzhe Liu, Xiaojuan Qi, and Philip H. S. Torr. Global texture enhancement for fake face detection in the wild. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 24

[47] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 7, 25

[48] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ODE solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 2022. 4, 16

[49] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 3, 16

[50] Zeyu Lu, Chengyue Wu, Xinyuan Chen, Yaohui Wang, Lei Bai, Yu Qiao, and Xihui Liu. Hierarchical diffusion autoencoders and disentangled image manipulation. *arXiv preprint arXiv:2304.11829*, 2023. 2

[51] Yanru Lyu, Xinxin Wang, Rungtai Lin, and Jun Wu. Communication in human–ai co-creation: Perceptual analysis of paintings generated by text-to-image system. *Applied Science*, 2022. 26, 29

[52] Francesco Marra, Diego Gragnaniello, Davide Cozzolino, and Luisa Verdoliva. Detection of gan-generated fake images over social networks. *IEEE Multimedia Information Processing and Retrieval*, 2018. 9, 26

[53] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 2022. 30

[54] Yisroel Mirsky and Wenke Lee. The creation and detection of deepfakes: A survey. *ACM Computing Surveys*, 2022. 26

[55] Lakshmanan Nataraj, Tajuddin Manhar Mohammed, B. S. Manjunath, Shivkumar Chandrasekaran, Arjuna Flenner, Jawadul H. Bappy, and Amit K. Roy-Chowdhury. Detecting GAN generated fake images using co-occurrence matrices. *Media Watermarking, Security, and Forensics*, 2019. 9, 26

[56] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *International Conference on Machine Learning*, 2021. 2

[57] Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. Towards universal fake image detectors that generalize across generative models. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 7, 9, 26, 30

[58] OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 25

[59] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2

[60] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in frequency: Face forgery detection by mining frequency-aware clues. *European Conference on Computer Vision*, 2020. 24

[61] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. *International Conference on Machine Learning*, 2021. 7, 25, 28

[62] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 2020. 25

[63] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *Advances in Neural Information Processing Systems*, 2022. 2, 9, 25, 26

[64] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 3, 9, 17, 25, 26, 30

[65] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arXiv:2208.12242*, 2022. 29

[66] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 2022. 2, 9, 25, 26

[67] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: an open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 2022. 9

[68] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. *Association for Computational Linguistics*, 2018. 4, 5, 18

[69] Philip L Smith and Daniel R Little. Small is beautiful: In defense of the small-n design. *Psychonomic bulletin & review*, 2018. 21

[70] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *International Conference on Learning Representations*, 2020. 2, 4, 16

[71] Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems*, 2021.

[72] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in Neural Information Processing Systems*, 2020.

[73] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *International Conference on Learning Representations*, 2021. 2, 25

[74] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *International Conference on Machine Learning*, 2021. 27

[75] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 25

[76] Luisa Verdoliva, Davide Cozzolino, and Koki Nagano. 2022 ieee image and video processing cup synthetic image detection. 2022. 5

[77] Ziyu Wan, Bo Zhang, Dongdong Chen, Pan Zhang, Dong Chen, Jing Liao, and Fang Wen. Bringing old photos back to life. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 29

[78] Run Wang, Felix Juefei-Xu, Lei Ma, Xiaofei Xie, Yihao Huang, Jian Wang, and Yang Liu. Fakespotter: A simple yet robust baseline for spotting ai-synthesized fake faces. *International Joint Conference on Artificial Intelligence*, 2020. 5

[79] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A. Efros. Cnn-generated images are surprisingly easy to spot... for now. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 5, 7, 9, 26, 30

[80] Yabin Wang, Zhiwu Huang, and Xiaopeng Hong. Benchmarking deepart detection. *arXiv preprint arXiv:2302.14475*, 2023. 5

[81] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *Advances in Neural Information Processing Systems*, 2022. 9, 25, 30

[82] Han Zhang, Jing Yu Koh, Jason Baldridge, Honglak Lee, and Yinfei Yang. Cross-modal contrastive learning for text-to-image generation. *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 2, 9, 25, 26

[83] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 2

[84] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. Detecting and simulating artifacts in GAN fake images. *IEEE International Workshop on Information Forensics and Security*, 2019. 9, 24, 26

[85] Yifan Zhang, Bingyi Kang, Bryan Hooi, Shuicheng Yan, and Jiashi Feng. Deep long-tailed learning: A survey. *arXiv preprint arXiv:2110.04596*, 2021. 30

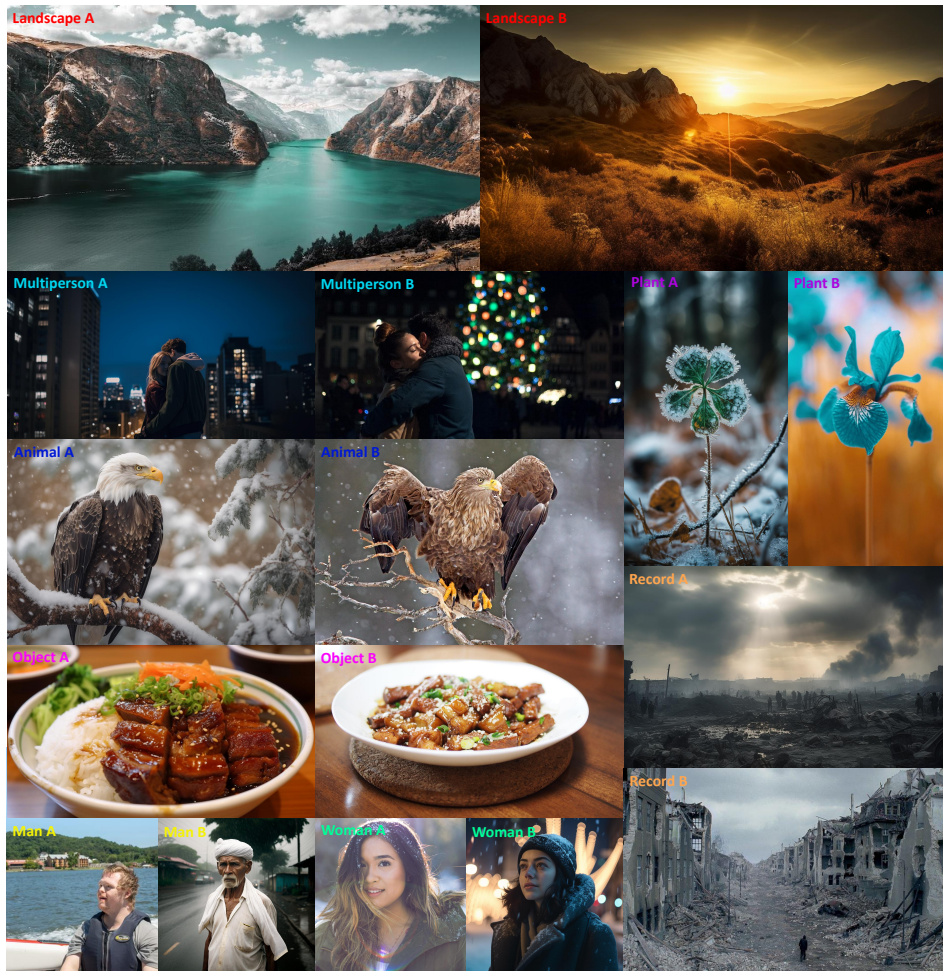## A    Quick Test: Can you identify which ones are AI-generated images?



Figure 5: **A quick test:** *Can you identify which ones are AI-generated images?*

**Answer of the Quick Test**    The AI-generated images of the Fig. 5 are "Landscape B", "Multiperson A", "Plant A", "Animal A", "Record A", "Object A", "Man B", "Woman B", respectively.

Table 8: **Detailed information of the datasets used in MPBench.** <span style="color:red">**R**</span> denotes the dataset consisting entirely of real images. <span style="color:red">**F**</span> denotes the dataset consisting entirely of fake images. ✔ denotes existing datasets. ✗ denotes the datasets provided in this work. "Diff" refers to diffusion model, "AR" refers to autoregressive model and "Unk." refers to unknown model. "Resolution" refers to the resolution of the fake images in the dataset. "Caption" refers to the caption used in text-to-image generation models to generate the corresponding dataset.

| Dataset | CC3M-Train | StyleGAN3-Train | SD-V1.5Real-dpms-25 | IF-V1.0-dpms++25 | StyleGAN3 | ImageNet-Test | CelebA-HQ-Train | CC3M-Val | SD-V2.1-dpms-25 | SD-V1.5-dpms-25 | SD-V1.5Real-dpms-25 | IF-V1.0-dpms++-10 | IF-V1.0-dpms++-25 | IF-V1.0-dpms++-50 | IF-V1.0-ddim-50 | IF-V1.0-dpms-50 | Cogview2 | Midjourney | StyleGan3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Category | Train | | | | | Validate | | | | | | | | | | | | | |
| | R | R | F | F | F | R | R | R | F | F | F | F | F | F | F | F | F | F | F |
| Generator | - | - | Diff. | Diff. | GAN | - | - | - | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | AR | Unk. | GAN |
| Numbers | 1M | 87K | 1M | 1M | 87K | 100K | 24K | 15K | 15K | 15K | 15K | 15K | 15K | 15K | 15K | 15K | 22K | 5.5K | 60K |
| Resolution | - | - | 512 | 256 | (>=)512 | - | - | - | 512 | 512 | 512 | 256 | 256 | 256 | 256 | 256 | 480 | (>=)640 | (>=)512 |
| Caption | - | - | CC3M-train(first 1M) | CC3M-train(first 1M) | - | - | - | - | CC3M-val | CC3M-val | CC3M-val | CC3M-val | CC3M-val | CC3M-val | CC3M-val | CC3M-val | - | - | - |
| Seed | - | - | 420 | 420 | - | - | - | - | 420 | 420 | 420 | 420 | 420 | 420 | 420 | 420 | - | - | - |
| CFG-Scale | - | - | 7 | 7 | - | - | - | - | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | - | - | - |
| This work | ✗ | ✗ | ✔ | ✔ | ✔ | ✗ | ✗ | ✗ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |

Table 9: **Detailed information of the diffusion datasets used in MPBench.**

| Category | Train | | Test | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Generators | SD-V1.5Real-dpms-25 | IF-V1.0-dpms++-25 | SD-V2.1-dpms-25 | SD-V1.5-dpms-25 | SD-V1.5Real-dpms-25 | IF-V1.0-dpms++-10 | IF-V1.0-dpms++-25 | IF-V1.0-dpms++-50 | IF-V1.0-ddim-50 | IF-V1.0-dpms-50 |
| Total Numbers | 1M | 1M | 15K | 15K | 15K | 15K | 15K | 15K | 15K | 15K |
| Sampling Steps | 25 | 25 | 25 | 25 | 25 | 10 | 25 | 50 | 50 | 50 |
| Sampling Methods | dpm-solver [48] | dpm-solver++ [49] | dpm-solver [48] | dpm-solver [48] | dpm-solver [48] | dpm-solver++ [49] | dpm-solver++ [49] | dpm-solver++ [49] | ddim [70] | ddpm [36]s |
| Seed | 420 | 420 | 420 | 420 | 420 | 420 | 420 | 420 | 420 | 420 |
| CFG-Scale | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |
| Model | Stable Diffusion v1.5 Realistic Version | IF v1.0 | Stable Diffusion v2.1 | Stable Diffusion v1.5 | Stable Diffusion v1.5 Realistic Version | IF v1.0 | IF v1.0 | IF v1.0 | IF v1.0 | IF v1.0 |

Table 10: **Detailed information of the StyleGAN3 datasets used in MPBench.**

| Category | Train | | | | | | Validate | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Generators | stylegan3-r-ffhqu-1024x1024 | stylegan3-t-ffhqu-1024x1024 | stylegan3-r-afhqv2-512x512 | stylegan3-t-afhqv2-512x512 | stylegan3-r-metfaces-1024x1024 | stylegan3-t-metfaces-1024x1024 | stylegan3-r-ffhqu-1024x1024 | stylegan3-t-ffhqu-1024x1024 | stylegan3-r-afhqv2-512x512 | stylegan3-t-afhqv2-512x512 | stylegan3-r-metfaces-1024x1024 | stylegan3-t-metfaces-1024x1024 |
| Total Numbers | 87K | | | | | | 60K | | | | | |
| Numbers | 35K | 35K | 8K | 8K | 0.65K | 0.65K | 10K | 10K | 10K | 10K | 10K | 10K |
| Seeds | 10001~45000 | 10001~45000 | 10001~18000 | 10001~18000 | 10001~10800 | 10001~10800 | 1-10000 | 1-10000 | 1-10000 | 1-10000 | 1-10000 | 1-10000 |
| Matched Dataset | FFHQ (70K) [42] | | AFHQv2 (16K) [19] | | Metfaces (1.3K) [41] | | None | | | | | |

## B    Dataset

### B.1    Dataset Configuration for Model Evaluation

We detailed the collection process of our datasets in Section 2.2 of the main paper, now the following will provide more detailed configuration information for each dataset.

We use the default Github repository code of each model to generate our datasets. Detailed information about the training and validation datasets are shown in Tab. 8. We further provide the captions and resolutions used in each specific dataset. For diffusion generation, we use the fixed seed and cfg-scale to generate our datasets. We also use different sampling methods and steps for generation. The detailed information about sampling methods and steps for different diffusion models can be found in Tab. 9. For StyleGAN3 generation, we use 2 models (stylegan3-r-ffhqu-1024x1024, stylegan3-t-ffhqu-1024x1024) to generate 70K face images for training to match the number of FFHQ and 10K face images for testing. We use 2 models (stylegan3-r-afhqv2-512x512, stylegan3-t-afhqv2-512x512) to generate 16K animal faces to match the number of AFHQ-v2 and 10K animal faces for testing. We use 2 models (stylegan3-r-metfaces-1024x1024, stylegan3-t-metfaces-1024x1024) to generate 1.3K art human faces for training to match the number of MetFaces Dataset and 10K art human faces for testing. The detailed information about our StyleGAN3 generation can be found in Tab. 10.

## B.2 Data Content Component Analysis of Training and Validation Dataset

We will analyze the composition of the training and validation dataset in the following two parts and discuss the issue of data imbalance. We also provide a detailed table showing the composition and proportion of different datasets, as shown in Tab. 11.

**Training Dataset.**

• **Fake2M** is composed of 1M fake images generated by the first 1M caption in CC3M using SD-V1.5Real-dpms-25 [64], 1M fake images generated by the first 1M caption in CC3M using IF-V1.0-dpms++-25 [5] and 87K fake images generated using StyleGAN3 [41], as shown in Tab. 8 and Tab 9.

In Fake2M, the number of face data is only 82K, accounting for %4 of the total data 2M, as shown in Tab. 11. There is no content imbalanced problem in Fake2M.

• **Training Dataset Setting A** is composed of 1M fake images generated by the first 1M caption in CC3M using SD-V1.5Real-dpms-25 in Fake2M and the first 1M real images in CC3M.

In Training Dataset Setting A, most of the content is general content. There is no content imbalanced problem in Training Dataset Setting A.

• **Training Dataset Setting B** is composed of 1M fake images generated by the first 1M caption in CC3M using IF-V1.0-dpms++-25 in Fake2M and the first 1M real images in CC3M.

In Training Dataset Setting B, most of the content is general content. There is no content imbalanced problem in Training Dataset Setting B.

• **Training Dataset Setting C** is composed of 87K fake images generated by StyleGAN3 in Fake2M (the detailed content in this dataset can be found in Tab. 10) and the first 1M real images in CC3M.

In training dataset setting C, most of the content is face. There is content imbalanced problem in training dataset setting C. This inclusion was intentional, aiming to specifically investigate the performance implications of face fake images produced by StyleGAN3.

• **Training Dataset Setting D** is composed of 460K fake images generated by the first 460K caption in CC3M using IF-V1.0-dpms++-25 in Fake2M, 460K fake images generated by the first 460K caption in CC3M using SD-V1.5Real-dpms-25 in Fake2M, 87K fake images generated by StyleGAN3, the first 1M real images in CC3M and 87K real images in StyleGAN3 training dataset.

In training dataset setting D, most of the content is general content. The number of fake face data is 82K and real face data is also 82K, accounting for %8 of the total data 2M. There is no content imbalanced problem in training dataset setting D.

**Validation Dataset (MPBench).** In MPBench, most of the content is general content. The number of fake face data is 60K, accounting for %15.3 of the total data 391.5K. The number of real face data is 24K, accounting for %6.1 of the total data 391.5K. There is no content imbalanced problem in training dataset setting D.

From the perspective of the ratio between real and fake images, we observe that the proportion of real and fake images is essentially the same across the four dataset settings and MPBench, as shown in Tab. 11. Therefore, there is no imbalance issue between the number of fake and real images.

## B.3 Quality Analysis

We conducted further analysis of our dataset quality score distributions, as shown in Fig. 6 and Tab. 12. We observed that the majority of our sub dataset have an average score above 0.6 (with Midjourneyv5-5K having an average score of 0.66) and the average score of all images in the dataset is 0.6. These demonstrate that our dataset is a high-quality dataset with a large amount of high-quality images. Only a few datasets (cogview2-22K, IF-ddim-25-15K-1024x1024, IF-ddim-50-15K-1024x1024, stylegan3-r-ffhqu-1024x1024, and stylegan3-r-metfaces-1024x1024) have an average score below 0.6. The distribution of quality scores across the entire dataset demonstrates a balanced mixture of high-quality and low-quality images, as shown in the "all-images" violin plot of Fig. 6. This aligns with our original intention: a fake image detection dataset should encompass both high-quality and low-quality image data. In order to better showcase our dataset, we provided more visualizations about the high quality, mid quality and low quality images in our dataset, as shown in Fig. 7.

Table 11: **Data Content Component Analysis.** "Content" means the type of the content in each dataset ("Face" means that the content in this dataset is mostly faces, such as FFHQ [42]. "Object" means that the content in this dataset is mainly composed of a limited number of objects, such as ImageNet [22]. "General" means that the content in this dataset is general, not limited to some objects, faces or art, such as CC3M [68]). "Each Dataset / Total Number (%)" means the number of images in this dataset and the percentage it contributes to the entire dataset setting setting. "Fake / Total Number (%)" means the number of fake images in the whole dataset setting and the percentage it contributes to the entire dataset setting. "Real / Total Number (%)" means the number of real images in the whole dataset setting and the percentage it contributes to the entire dataset setting.

| Dataset | Fake | | | | Real | | | |
|---|---|---|---|---|---|---|---|---|
| | Name | Content | Each Dataset / Total Number (%) | Fake / Total Number (%) | Name | Content | Each Dataset / Total Number (%) | Real / Total Number (%) |
| Fake2M Dataset | SD-V1.5Real-dpms-25 | General | 1M (47.9%) | | | | | |
| | IF-V1.0-dpms++-25 | General | 1M (47.9%) | 2.08M (100%) | | | | |
| | StyleGAN3 | Face | 87K (4.2%) | | | | | |
| Training Dataset Setting A | SD-V1.5Real-dpms-25 | General | 1M (50%) | 1M (50%) | CC3M-Train | General | 1M (50%) | 1M (50%) |
| Training Dataset Setting B | IF-V1.0-dpms++-25 | General | 1M (50%) | 1M (50%) | CC3M-Train | General | 1M (50%) | 1M (50%) |
| Training Dataset Setting C | StyleGAN3 | Face | 87K (50%) | 87K (50%) | CC3M-Train | General | 87K (50%) | 87K (50%) |
| Training Dataset Setting D | SD-V1.5Real-dpms-25 | General | 460K (21.2%) | | CC3M-Train | General | 1M (46%) | |
| | IF-V1.0-dpms++-25 | General | 460K (21.2%) | 1.08M (50%) | StyleGAN3-Train | Face | 87K (4%) | 1.08M (50%) |
| | StyleGAN3 | Face | 87K (4%) | | | | | |
| Validation Dataset (MPBench) | SD-V2.1-dpm-25 | General | 15K (3.8%) | | ImageNet-Test | Object | 100K (25.5%) | |
| | SD-V1.5-dpm-25 | General | 15K (3.8%) | | CelebA-HQ-Train | Face | 24K (6.1%) | |
| | SD-V1.5Real-dpm-25 | General | 15K (3.8%) | | CC3M-Val | General | 15K (3.8%) | |
| | IF-V1.0-dpm++-10 | General | 15K (3.8%) | | | | | |
| | IF-V1.0-dpm++-25 | General | 15K (3.8%) | | | | | |
| | IF-V1.0-dpm++-50 | General | 15K (3.8%) | 252.5K (64.5%) | | | | 139K (35.5%) |
| | IF-V1.0-ddim-50 | General | 15K (3.8%) | | | | | |
| | IF-V1.0-ddpm-50 | General | 15K (3.8%) | | | | | |
| | Cogview2 | General | 22K (5.6%) | | | | | |
| | Midjourney | General | 5.5K (1.4%) | | | | | |
| | StyleGAN3 | Face | 60K (15.3%) | | | | | |

Table 12: **Quality score distribution statistical information of the dataset.** "all-images" means the quality score distribution of all the images in the dataset. "Mean Score" means the average score of the quality score in the sub dataset. "Min Score" means the minimum score of the quality score in the sub dataset. "Max Score" means the maximum score of the quality score in the sub dataset.

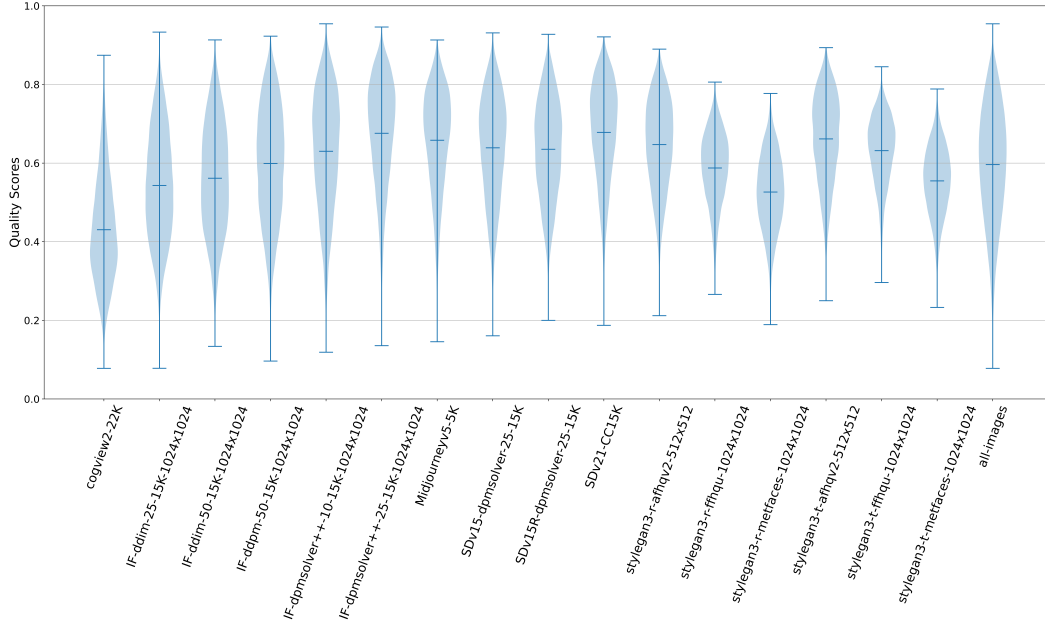| Sub Dataset | Mean Score | Min Score | Max Score |
|---|---|---|---|
| cogview2-22K | 0.43 | 0.08 | 0.87 |
| IF-ddim-25-15K-1024x1024 | 0.54 | 0.08 | 0.93 |
| IF-ddim-50-15K-1024x1024 | 0.56 | 0.13 | 0.91 |
| IF-ddpm-50-15K-1024x1024 | 0.60 | 0.10 | 0.92 |
| IF-dpmsolver++-10-15K-1024x1024 | 0.63 | 0.12 | **0.95** |
| IF-dpmsolver++-25-15K-1024x1024 | 0.68 | 0.14 | **0.95** |
| Midjourneyv5-5K | 0.66 | 0.15 | 0.91 |
| SDv15-dpmsolver-25-15K | 0.64 | 0.16 | 0.93 |
| SDv15R-dpmsolver-25-15K | 0.64 | 0.20 | 0.93 |
| SDv21-CC15K | **0.68** | 0.19 | 0.92 |
| stylegan3-r-afhqv2-512x512 | 0.65 | 0.21 | 0.89 |
| stylegan3-r-ffhqu-1024x1024 | 0.59 | 0.27 | 0.81 |
| stylegan3-r-metfaces-1024x1024 | 0.53 | 0.19 | 0.78 |
| stylegan3-t-afhqv2-512x512 | 0.66 | 0.25 | 0.89 |
| stylegan3-t-ffhqu-1024x1024 | 0.63 | **0.30** | 0.85 |
| stylegan3-t-metfaces-1024x1024 | 0.55 | 0.23 | 0.79 |
| all-images | 0.60 | 0.08 | **0.95** |

18

Figure 6: **Quality score distribution of the dataset.** "all-images" means the quality score distribution of all the images in the dataset.



| IQA score | 0.396 0.673 0.868 | 0.313 0.680 0.929 | 0.370 0.640 0.901 |
| | cogview2-22K | IF-dpmsovler++-10-15K | SDv15-dpmsovler-25-15K |

| IQA score | 0.355 0.694 0.896 | 0.382 0.680 0.938 | 0.353 0.605 0.902 |
| | IF-ddim-25-15K | IF-dpmsovler++-25-15K | SDv15R-dpmsovler-25-15K |

| IQA score | 0.346 0.623 0.896 | 0.381 0.653 0.931 | 0.352 0.635 0.918 |
| | IF-ddim-50-15K | IF-dpmsovler++-50-15K | SDv21-dpmsovler-25-15K |

| IQA score | 0.376 0.689 0.909 | 0.386 0.678 0.900 | 0.378 0.664 0.889 |
| | IF-ddpm-50-15K | Midjourney-5.5K | Stylegan3-60K |

Figure 7: **Image visualization with image quality score.**

## C  HPBench

### C.1  Procedures for HPBench

Fig. 8 shows that our HPBench could be divided into three parts. In the first part, we collect realistic AI-generated images and real images across eight categories using the expertise of an annotator to

filter out low-quality AI-generated images. In the second part, we recruit a total of 50 volunteers for human evaluation. For each volunteer, he/she should complete a 100-question questionnaire in our prepared environment. In the third part, We disposal and analyze human evaluation data to draw conclusions.
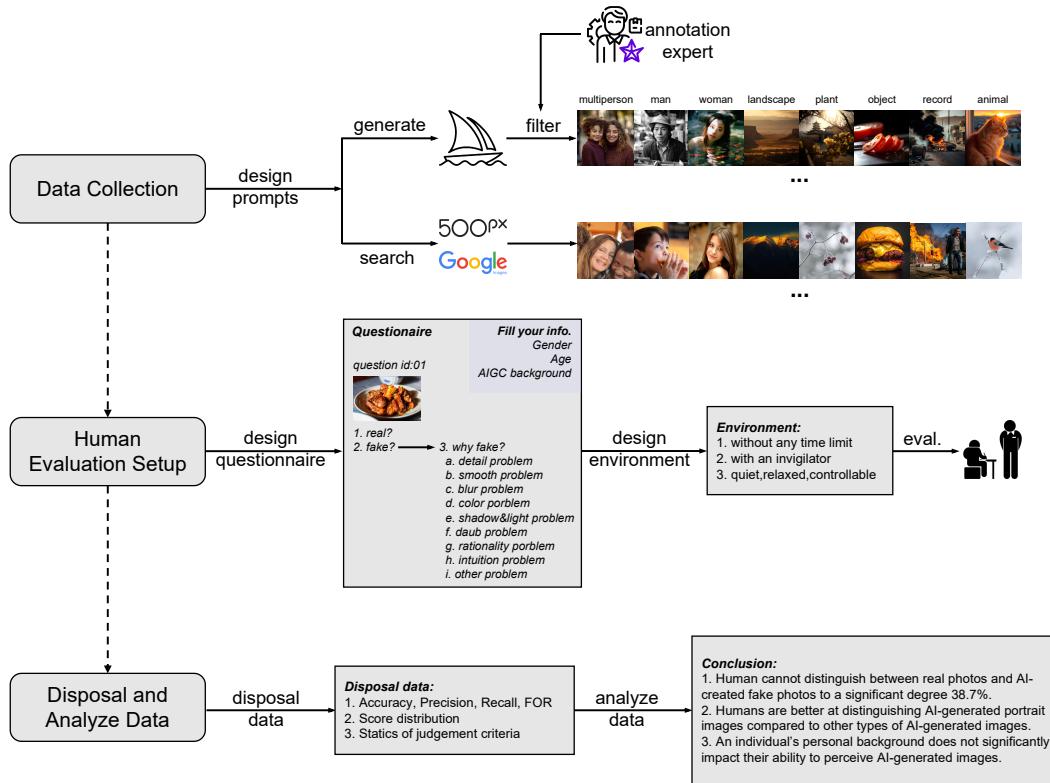


Figure 8: **Procedures for HPBench.**

## C.2   Detail of Human Evaluation

We recruit a total of 50 participants to participate in our human evaluation instead of using crowd-sourcing. In order to ensure comprehensiveness, fairness, and quality of the evaluation, we make efforts to ensure the diversity of the participants. Each participant is asked to complete a questionnaire consisting of 100 questions without any time limit. The questionnaires are completed in the presence of a project team member to guide the participant and ensure the quality of the human evaluation results. It is worth noting that we did not inform the participant about the ratio of real photos to AI-generated images in the questionnaires.

Specifically, each question in the questionnaire provides the participant with an image, and the participant is asked to determine whether the image is generated by AI or not. If the participant thinks that the image is generated by AI, he/she will be required to choose one or more reasons from the eight predefined judgment criteria options or provide their own judgment criteria. The eight options are explained below:

• Detail: AI-generated images may lack fine details, such as wrinkles in clothing or hair details.

• Smooth: AI-generated images may appear smoother or more uniform than real photos, such as smooth skin or unrealistic facial expressions.

• Blur: AI-generated images may be blurry, such as blurry or unclear edges.

• Color: AI-generated images may have unrealistic or inconsistent colors, such as colors that are too bright, too dark, or like the color of animation.

- Shadow & Light: AI-generated images may have unrealistic or inconsistent shadows and lighting, such as shadows or lightning that violate physics.

- Daub: AI-generated images may contain rough, uneven, or poorly applied colors or textures.

- Rationality: AI-generated images may contain irrational/illogical/contradictory contents.

- Intuition: people may judge whether a photo is AI-generated or not by intuition and cannot describe the exact reasons.

It is important to point out that each questionnaire consists of 50 real images and 50 AI-generated images, all of which are randomly sampled from the real image database containing 244 images and the AI-generated image database containing 151 images.

### C.3 Reason of High-quality Fifty-participant Human Evaluation

Inspired by Robert *et al.* [30], we aim to collect high-quality human evaluation data instead of noisy crowdsourcing data to ensure the high quality of the results. Our human evaluation, with a limited number of participants in a controlled environment, and conducting multiple experiments, is commonly referred to as the "small-N design". As the article "Small is beautiful: In defense of the small-N design" [69] suggests, the "small-N design" is the core of high-quality psychophysics. Crowdsourcing involves many participants in an uncontrollable, noisy setting, with each performing fewer trials. In contrast, we strive to ensure that each participant, with diverse backgrounds, takes the full 100-question survey in a quiet, relaxed, controllable, and monitored environment. Those factors contribute to high data quality.

### C.4 Crowd Sourcing Human Evaluation

We collect 1085 crowd-sourced human evaluation questionnaires to make the entire benchmark more comprehensive. We utilize the same experimental setup as the "High-quality Fifty-participant Human Evaluation" before. As the questionnaires are obtained through crowd-sourcing in an uncontrollable and noisy setting, we do not ask the participants to provide justification for each decision. We collect the accuracy of all the questionnaires, and the accuracy is only 49.9%. This highlights that in a fast-paced, noisy, and uncontrolled environment, people are completely unable to distinguish high-quality AI-generated images.

### C.5 Metrics of HPBench

We employ four commonly used evaluation metrics to analyze our results and highlight their respective meanings in the context of our problem. We define positive samples as AI-generated images and negative samples as real images for our problem, and then calculate Accuracy, Precision, Recall, and False Omission Rate (FOR) in the context of our problem.

**Accuracy** is a statistical measure used to evaluate how well a binary classification test correctly identifies or excludes a condition. In our study, accuracy represents the average precision of humans in distinguishing AI-generated images from real images.

**Precision** is the percentage of predicted positive cases that are actually positive. In our study, high precision represents the proportion of AI-generated images out of the total number of images that are predicted as AI-generated ones.

**Recall** is the percentage of true positive cases that are actually predicted as positive. In our study, recall represents the proportion of AI-generated images that are correctly identified as such out of the total number of AI-generated images.

**FOR** is the percentage of false negatives out of all negative cases. In our problem, FOR represents the proportion of real images misidentified as AI-generated images out of the total number of images.

### C.6 Analysis of the AIGC defects

Based on the user data we collected above, we summarize and show nine shortcomings of the current AIGC, as shown in Fig. 9: (1) "Hand problem" refers to situations where fingers overlap or have unreasonable shapes (multi fingers), resulting in images that are not realistic. (2) "Smoothing
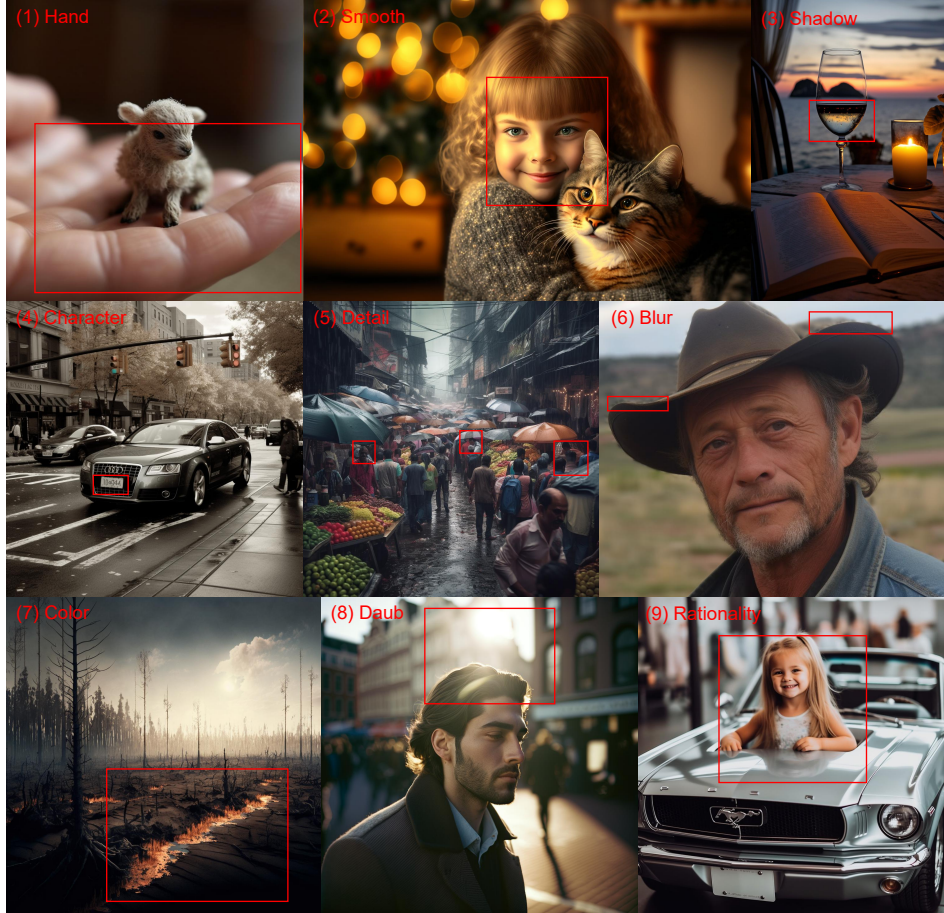
Figure 9: **Nine shortcomings of current AIGC.** We highlight the obvious defects of AIGC with red boxes.

problem" refers to situations where AI-generated images have overly smooth skin, resulting in unrealistic facial expressions and features. (3) "Shadow&Light problem" refers to situations where the position and shape of light sources and shadows in AI-generated images are unreasonable, resulting in unnatural lighting effects. (4) "Character problem" refers to situations where incorrect signs and texts appear in AI-generated images, which do not match reality. (5) "Detail problem" refers to situations where some details in AI-generated images are not realistic or unreasonable, such as wrinkles in clothing or hair details. (6) "Blur problem" refers to situations where AI-generated images are blurry or unclear, resulting in obvious artifacts. (7) "Color problem" refers to situations where the colors in AI-generated images are not realistic or coordinated, such as colors that are too bright, too dark, or like the color of animation. (8) "Daub problem" refers to situations where AI-generated images have been excessively daubed, resulting in lost details or unrealistic images. (9) "Rationality problem" between objects refers to situations where the relationship between objects in AI-generated images is not reasonable, such as incorrect size proportions or unreasonable positions.

## C.7 Detailed Score Distribution

**Detailed score distribution of different categories for all volunteers.** As shown in Fig. 12, we visualize the detailed score distribution of different categories for all volunteers: (a) the detailed score distribution of different categories for all volunteers and all tested images (b) the detailed score distribution of different categories for all volunteers and only AI-generated images (c) the detailed score distribution of different categories for all volunteers and only real images.

**Detailed score distribution of different categories for men and women.** As shown in Fig. 10, we visualize the score distribution of all images for man and woman. We find that the average scores of men and women are almost the same with a relatively accuracy rate of 61%.

We also visualize the detailed score distribution of different categories for man and woman in Fig. 13: (a) the detailed score distribution of different categories for man and woman and all tested images (b) the detailed score distribution of different categories for man and woman and only AI-generated images (c) the detailed score distribution of different categories for man and woman and only real images. A interesting finding is that: Apart from humans having a high recognition rate for human portraits, men have a higher recognition rate for the category *Man* than women, and women have a higher recognition rate for the category *Woman* than men. We speculate that people may have a higher recognition accuracy for more familiar objects.
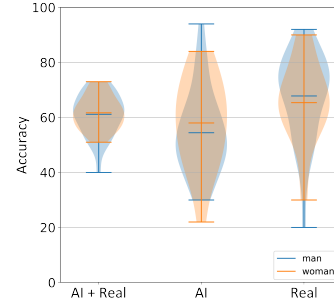


Figure 10: **Human evaluation score distributions calculated by data from men and women.**

**Detailed score distribution of different categories for volunteers with and without AIGC background.** As shown in Fig. 14, we visualize the detailed score distribution of different categories for volunteers with and without AIGC background: (a) the detailed score distribution of different categories for volunteers with and without AIGC background and all tested images (b) the detailed score distribution of different categories for volunteers with and without AIGC background and only AI-generated images (c) the detailed score distribution of different categories for volunteers with and without AIGC backgrounds and only real images.

# D  MPBench

## D.1  More Experiments on MPBench

We conducted more experiments on MPBench, as shown in Tab. 13.

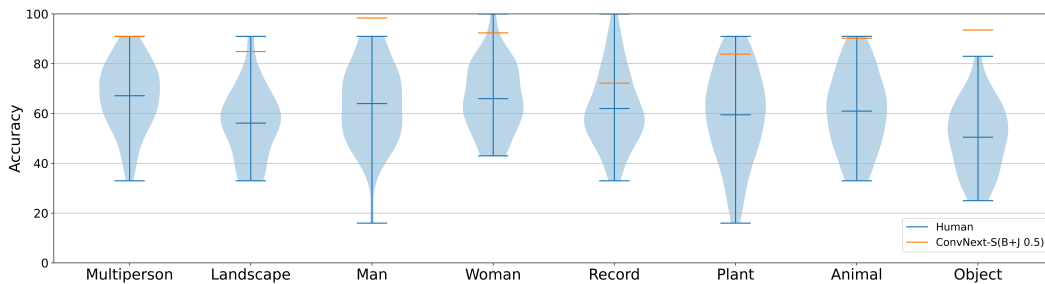## D.2  Evaluate the best model under the same setting used in HPBench.



Figure 11: **Human evaluation score distribution and ConvNext-S(B+J 0.5) model score in the same dataset HPBench.**

We also present the results of human and ConvNext-S(B+J 0.5) model on the same dataset HPBench, as shown in Fig. 11. We can find that the results of ConvNext-S are better than human in the results of the two categories: *Man* and *Object*. In the remaining categories, the highest performance of human is better than the performance of the model, but the average performance of human is far worse than the performance of the model.

This demonstrates that it is valuable to study the potential benefits of ensemble the abilities of humans and models in addressing this challenge.

Table 13: **Quantitative comparison of another five models under four training dataset settings with fourteen validation datasets.** "Diff" refers to diffusion model, "AR" refers to autoregressive model and "Unk." refers to unknown model. **Real (R)** denotes the dataset consisting entirely of real images. **Fake (F)** denotes the dataset consisting entirely of fake images. Blue cell denotes the deepfake methods.

| Model | Training Dataset | ImageNet-Test | CelebA-HQ-Train | CC3M-Val | Average Acc. | SD-V2.1-dpm-25 | SD-V1.5-dpm-25 | SD-V1.5Real-dpm-25 | IF-V1.0-dpm+-10 | IF-V1.0-dpm+-25 | IF-V1.0-dpm++-50 | IF-V1.0-ddim-50 | IF-V1.0-ddpm-50 | Cogview2 | Midjourney | StyleGAN3 | Average Acc. | Total Average Acc. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | **Real** | | | | **Fake** | | | | | | | | | | | | **R+F** |
| | | - | - | - | - | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | Diff. | AR | Unk. | GAN | - | - |
| Swin-S(B+J 0.1) | | 92.4 | 100.0 | 99.9 | 97.4 | 37.9 | 71.7 | 99.9 | 94.4 | 69.6 | 69.4 | 50.6 | 76.9 | 49.5 | 6.8 | 39.5 | 60.5 | 68.4 |
| Swin-S(B+J 0.5) | | 95.1 | 99.9 | 99.9 | 98.3 | 41.3 | 60.2 | 99.9 | 89.7 | 52.3 | 52.2 | 42.0 | 67.6 | 64.5 | 9.5 | 25.1 | 54.9 | 64.2 |
| DeiT-S(B+J 0.1) | Dataset Setting A: | 99.7 | 99.9 | 99.9 | 99.8 | 37.8 | 47.3 | 99.9 | 19.6 | 2.6 | 2.2 | 5.2 | 12.9 | 8.5 | 6.2 | 2.1 | 22.2 | 38.8 |
| DeiT-S(B+J 0.5) | SD-V1.5Real-dpms-25 (1M) CC3M-Train (1M) | 99.4 | 99.9 | 99.8 | 99.7 | 46.2 | 51.5 | 99.9 | 21.5 | 4.0 | 3.1 | 6.6 | 13.9 | 4.7 | 8.6 | 2.8 | 23.8 | 40.1 |
| ResNet50(Fourier) [84] | | 42.3 | 7.2 | 4.3 | 17.9 | 95.9 | 97.4 | 97.3 | 96.9 | 96.0 | 95.4 | 98.6 | 96.0 | 95.6 | 94.9 | 97.0 | 96.4 | 79.6 |
| Xception(Patch) [16] | | 57.1 | 58.7 | 55.7 | 57.1 | 31.0 | 34.9 | 31.4 | 36.5 | 34.7 | 36.9 | 35.5 | 37.6 | 49.7 | 39.5 | 42.7 | 37.3 | 41.5 |
| Swin-S(B+J 0.1) | | 88.1 | 99.9 | 99.9 | 95.9 | 0.4 | 3.0 | 0.1 | 99.6 | 99.8 | 99.7 | 2.1 | 22.0 | 3.3 | 0.1 | 0.7 | 30.0 | 44.1 |
| Swin-S(B+J 0.5) | | 81.5 | 99.9 | 99.9 | 93.7 | 0.9 | 9.9 | 0.3 | 99.6 | 99.9 | 99.8 | 6.7 | 42.6 | 6.8 | 0.1 | 3.3 | 33.6 | 46.5 |
| DeiT-S(B+J 0.1) | Dataset Setting B: | 98.3 | 99.7 | 99.7 | 99.2 | 8.9 | 32.0 | 9.5 | 95.3 | 99.2 | 97.4 | 53.6 | 55.0 | 25.4 | 2.8 | 2.5 | 43.7 | 55.6 |
| DeiT-S(B+J 0.5) | IF-V1.0-dpms++-25 (1M) CC3M-Train (1M) | 97.8 | 99.6 | 99.4 | 98.9 | 11.3 | 36.6 | 10.0 | 95.2 | 99.3 | 97.9 | 57.9 | 61.9 | 27.8 | 3.7 | 2.9 | 45.8 | 57.2 |
| ResNet50(Fourier) [84] | | 42.3 | 45.7 | 51.0 | 46.3 | 60.7 | 61.7 | 73.0 | 59.3 | 72.2 | 70.3 | 29.5 | 69.4 | 40.8 | 60.9 | 70.5 | 60.7 | 57.6 |
| Xception(Patch) [16] | | 54.5 | 17.0 | 29.0 | 33.5 | 43.4 | 44.1 | 49.8 | 23.7 | 19.3 | 20.1 | 63.5 | 35.8 | 49.7 | 57.0 | 58.1 | 42.2 | 40.3 |
| Swin-S(B+J 0.1) | | 99.6 | 99.9 | 99.5 | 99.6 | 2.1 | 3.0 | 1.7 | 1.3 | 2.8 | 2.0 | 6.1 | 5.3 | 27.4 | 1.6 | 99.2 | 13.8 | 32.2 |
| Swin-S(B+J 0.5) | | 99.6 | 99.9 | 99.1 | 99.5 | 1.6 | 2.7 | 1.6 | 2.6 | 4.3 | 3.8 | 4.5 | 3.3 | 23.4 | 1.7 | 99.3 | 13.5 | 31.9 |
| DeiT-S(B+J 0.1) | Dataset Setting C: | 97.7 | 99.8 | 89.9 | 95.8 | 10.5 | 13.4 | 11.2 | 12.4 | 17.8 | 15.4 | 14.9 | 14.5 | 36.6 | 15.9 | 95.9 | 23.5 | 38.9 |
| DeiT-S(B+J 0.5) | StyleGAN3 (87K) StyleGAN3-Train (87K) | 97.1 | 99.5 | 87.1 | 94.5 | 11.3 | 13.6 | 10.4 | 13.4 | 20.7 | 17.6 | 15.8 | 15.2 | 39.5 | 18.5 | 94.6 | 24.6 | 39.5 |
| ResNet50(Fourier) [84] | | 99.2 | 99.9 | 99.2 | 99.4 | 0.1 | 0.2 | 0.1 | 0.6 | 0.2 | 0.1 | 0.2 | 0.5 | 1.0 | 0.6 | 0.45 | 0.3 | 21.5 |
| Xception(Patch) [16] | | 53.3 | 51.8 | 51.2 | 52.1 | 49.7 | 48.7 | 48.5 | 48.6 | 51.2 | 51.6 | 45.1 | 50.8 | 54.7 | 45.9 | 49.0 | 49.4 | 50.0 |
| F³-Net [60] | | 91.7 | 98.8 | 86.2 | 92.2 | 15.4 | 13.3 | 9.5 | 9.7 | 17.6 | 15.3 | 19.9 | 21.7 | 29.4 | 18.1 | 97.5 | 24.3 | 38.8 |
| Gramnet [46] | | 85.1 | 99.5 | 82.2 | 88.9 | 14.3 | 14.6 | 9.5 | 13.1 | 22.0 | 21.0 | 19.2 | 18.4 | 39.9 | 27.1 | 94.3 | 26.6 | 40.0 |
| ELA-Xception [32] | | 73.8 | 99.7 | 68.6 | 80.7 | 49.1 | 38.2 | 35.1 | 41.4 | 42.1 | 41.1 | 47.3 | 53.1 | 73.2 | 38.5 | 89.3 | 49.8 | 56.4 |
| Swin-S(B+J 0.1) | Dataset Setting D: | 83.2 | 99.9 | 99.9 | 94.3 | 48.9 | 92.3 | 99.9 | 99.9 | 99.9 | 99.9 | 67.9 | 95.1 | 61.4 | 13.3 | 99.3 | 79.8 | 82.9 |
| Swin-S(B+J 0.5) | SD-V1.5Real-dpms-25 (460K) | 93.5 | 99.9 | 99.9 | 97.7 | 47.2 | 79.9 | 99.9 | 99.7 | 99.8 | 99.7 | 59.1 | 93.5 | 64.1 | 10.6 | 98.8 | 77.4 | 81.8 |
| DeiT-S(B+J 0.1) | IF-V1.0-dpms++-25 (460K) | 98.2 | 99.9 | 99.6 | 99.2 | 51.0 | 69.4 | 99.8 | 94.6 | 98.4 | 95.0 | 48.0 | 56.5 | 37.3 | 10.3 | 96.7 | 68.6 | 75.3 |
| DeiT-S(B+J 0.5) | StyleGAN3 (87K) CC3M-Train (1M) | 96.7 | 99.7 | 99.0 | 98.4 | 54.8 | 75.7 | 99.8 | 96.0 | 99.1 | 97.6 | 65.7 | 79.1 | 55.4 | 12.1 | 93.1 | 75.3 | 80.2 |
| ResNet50(Fourier) [84] | StyleGAN3-Train (87K) | 58.0 | 26.8 | 62.5 | 49.1 | 41.2 | 51.2 | 48.4 | 43.2 | 50.7 | 53.8 | 7.3 | 53.1 | 21.4 | 58.7 | 64.3 | 44.8 | 45.7 |
| Xception(Patch) [16] | | 75.3 | 69.3 | 69.8 | 71.4 | 51.9 | 54.7 | 50.3 | 55.5 | 59.0 | 58.2 | 32.3 | 55.3 | 47.4 | 43.9 | 17.0 | 47.7 | 52.8 |

## D.3 Hyperparameters of the Experiments

Detailed information about the hyperparameters of the experiments in MPBench are shown in Tab. 14, Tab. 15, Tab. 16, Tab. 17 and Tab. 18.

| config | value |
|---|---|
| optimizer | AdamW |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ |
| weight decay | 0.05 |
| learning rate | 1e-4 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip, RandomResizedCrop(224), GaussianBlur(0.1), JPEG(0.1) |
| batch size | 1024 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | ConvNext-Small-In21k |

(a) ConvNext-S(B+J 0.1)

| config | value |
|---|---|
| optimizer | AdamW |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ |
| weight decay | 0.05 |
| learning rate | 1e-4 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip, RandomResizedCrop(224), GaussianBlur(0.5), JPEG(0.5) |
| batch size | 1024 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | ConvNext-Small-In21k |

(b) ConvNext-S(B+J 0.5)

Table 14: **Settings for ConvNext-S [47] in MPBench.**

# E    Detailed Related Work

## E.1    Image Generation

Generating photorealistic images based on given text descriptions has proven to be a challenging task. Previous GAN-based approaches [14, 31, 42, 82] were only effective within specific domains and datasets, assuming a closed-world setting. However, with the advancements in diffusion models [36, 73], autoregressive transformers [75], and large-scale language encoders [15, 58, 61, 62], significant progress has been made in high-quality photorealistic text-to-image synthesis with arbitrary text descriptions.

State-of-the-art text-to-image synthesis approaches such as DALL·E 2 [63], Imagen [66], Stable Diffusion [64], and Midjourney [7] have demonstrated the possibility of that generating high-quality, photorealistic images with diffusion-based generative models trained on large datasets. Those models have surpassed previous GAN-based models in both fidelity and diversity of generated images, without the instability and mode collapse issues that GANs are prone to. In addition to diffusion models, other autoregressive models such as Make-A-Scene [27], CogView [24], and Parti [81] have also achieved amazing performance. While diffusion models and autoregressive models exhibit impressive image synthesis ability, they all require time-consuming iterative processes to achieve high-quality image sampling. However, the progress made in the field of text-to-image synthesis over the past few years is a testament to the potential of this technology.

## E.2    Deepfake Generation and Detection

In December 2017, a Reddit user going by the pseudonym "Deepfakes" shared pornographic videos created using open-source AI tools capable of swapping faces in images and videos. Since then, the

| config | value |
|---|---|
| optimizer | AdamW |
| optimizer momentum | $\beta_1, \beta_2{=}0.9, 0.999$ |
| weight decay | 0.05 |
| learning rate | 1e-4 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip, RandomResizedCrop(224), GaussianBlur(0.1), JPEG(0.1) |
| batch size | 1024 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | Swin-Small-In1k |

(a) Swin-S(B+J 0.1)

| config | value |
|---|---|
| optimizer | AdamW |
| optimizer momentum | $\beta_1, \beta_2{=}0.9, 0.999$ |
| weight decay | 0.05 |
| learning rate | 1e-4 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip, RandomResizedCrop(224), GaussianBlur(0.5), JPEG(0.5) |
| batch size | 1024 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | Swin-Small-In1k |

(b) Swin-S(B+J 0.5)

Table 15: **Settings for Swin-S [45] in MPBench.**

term "Deepfake" has been widely used to describe the generation of human appearances, particularly facial expressions, through AI methods. The "Malicious Deep Fake Prohibition Act" of 2018 provides a definition of deepfake as videos and audios that have been realistically but falsely altered and are difficult to identify. Similarly, the "DEEP FAKES Accountability Act" of 2019 defines deepfake as media that is capable of authentically depicting an individual who did not actually participate in the production of the content. Yisroel *et al.* [54] defines deepfake as believable media generated by a deep neural network. In essence, deepfake [18] refers to the creation of seemingly realistic but falsified images, audios, videos, and other digital media produced through AI methods, particularly deep learning.

Realistic deepfake media has posed a significant threat to privacy, democracy, national security, and society as a whole. These images and videos have the potential to bypass facial authentication, create political unrest, spread fake news, and even be used for blackmail. The proliferation of fake information through fabricated videos and images can severely undermine our trust in online digital content. Furthermore, the highly realistic nature of deepfake media makes it difficult for humans to identify them as being falsified. Thus, the ability to distinguish between deepfake and real media has become an important, necessary, and urgent matter.

In recent years, there have been many works [11, 16, 20, 26, 52, 55, 57, 79, 84] exploring how to distinguish whether an image is AI-generated. These works focus on images generated by GANs or small generation models [14, 31, 42, 82]. Due to the limited quality of images generated by those methods, it is easy for humans to distinguish whether a photo is AI-generated or not. However, as the quality of generated images continues to improve with the advancement of recent generative models [7, 63, 64, 66], it has become increasingly difficult for humans to identify whether an image is generated by AI or not. Lyu *et al.* [51] provides an in-depth investigation into communication in human-AI co-creation, specifically focusing on the perceptual analysis of paintings generated by

| config | value |
|---|---|
| optimizer | AdamW |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ |
| weight decay | 0.05 |
| learning rate | 1e-4 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip, RandomResizedCrop(224), GaussianBlur(0.1), JPEG(0.1) |
| batch size | 1024 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | DeiT-Small-In1k |

(a) DeiT-S(B+J 0.1)

| config | value |
|---|---|
| optimizer | AdamW |
| optimizer momentum | $\beta_1, \beta_2 = 0.9, 0.999$ |
| weight decay | 0.05 |
| learning rate | 1e-4 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip, RandomResizedCrop(224), GaussianBlur(0.5), JPEG(0.5) |
| batch size | 1024 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | DeiT-Small-In1k |

(b) DeiT-S(B+J 0.5)

Table 16: **Settings for DeiT-S [74] in MPBench.**

a text-to-image system. Instead of exploring the human perception of AI-generated paintings, we study the human perception of AI-generated photographic images that may contain contradictions or absurdities that violate reality. Those AI-generated photorealistic images can potentially pose a significant threat to the accuracy of factual information. In conclusion, the objective of our study is to investigate whether state-of-the-art AI-generated photographic images are capable of deceiving human perception.

# F  Discussion, Broader Impact, Limitation and Conclusion

## F.1  Discussion

**Can AIGC deceive humans now?**   With the recent rapid advancements in generative AI, AI is now capable of producing highly photorealistic images with rich backgrounds, vivid characters, and beautiful lighting. Although people may able to occasionally differentiate low-quality AI-generated images, it is becoming more and more difficult to distinguish high-quality AI-generated images from real photography. In this study, our human evaluation results indicate that the state-of-the-art (SOTA) AI model is able to deceive the human eye to a significant degree (38.7%). Moreover, our exploration shows that it is no longer reliable to judge whether an image is real based solely on image quality. Instead, people need to consider factors such as over-smoothing portrait faces, coherence, and consistency between objects, and physical laws in the image, which makes the distinguishing process much harder and time-consuming (about 18 seconds for each image in this study).

From another aspect, current AI still can not **consistently** deceive the human eye. AI-generated images still have certain defeats which could be used by humans to distinguish fake images. Besides, creating such high-quality images requires prompt engineering skills and numerous experiments.

| config | value |
| --- | --- |
| optimizer | SGD |
| optimizer momentum | $\beta$=0.9 |
| weight decay | 1e-4 |
| learning rate | 1e-4 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip(0.5), RandomResizedCrop(224), GaussianBlur(0.1), JPEG(0.1) |
| batch size | 512 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | ResNet-50-In1k |

(a) ResNet50(B+J 0.1)

| config | value |
| --- | --- |
| optimizer | SGD |
| optimizer momentum | $\beta$=0.9 |
| weight decay | 1e-4 |
| learning rate | 1e-4 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip(0.5), RandomResizedCrop(224), GaussianBlur(0.5), JPEG(0.5) |
| batch size | 512 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | ResNet-50-In1k |

(b) ResNet50(B+J 0.5)

Table 17: **Settings for ResNet50 [33] in MPBench.**

| config | value |
| --- | --- |
| optimizer | AdamW |
| optimizer momentum | $\beta_1, \beta_2$=0.9, 0.999 |
| weight decay | 0.3 |
| learning rate | 1e-5 |
| learning rate sch. | cosine decay |
| warmup epochs | 0 |
| epochs | 10 |
| augmentation | HFlip(0.5), RandomResizedCrop(224) |
| batch size | 512 |
| dtype | bfloat16 |
| resolution | 224 |
| pretrain | openclip-ViT-L-14 |

Table 18: **Settings for CLIP-ViT-L [17, 61] in MPBench.**

Even though, a few finely adjusted AI images with misleading information can convey wrong ideas and cause enormous damage.

**What the current state-of-the-art image generation model can do and can not do?**　Given suitable prompts, the SOTA image generation model can produce photo-realistic images that are indistinguishable from real photographs, as shown in Fig. 5. The prompt can have different formats (e.g., text, image) and arbitrary complexity, including details such as colors, textures, and lighting. There are lots of potential applications for image generation. For instance, AIGC can be used to generate images for advertising campaigns, product catalogs, and fashion magazines. Since it can

easily be controlled by text, AIGC can also be utilized in the film industry to create realistic special effects or even entire scenes, at an extremely low cost. Furthermore, AIGC can be implemented in the gaming industry to produce immersive and lifelike game worlds.

Although generative AI has impressive image generation capabilities, it currently faces several limitations and challenges, as shown in Fig. 9. One of the most significant challenges is generating images of multiple people with intricate details in a single scene. Users can easily infer the authenticity of an image from details. Furthermore, the current model has difficulty generating realistic human hand gestures and positions, which are crucial for many applications such as sign language recognition and virtual reality. In addition, the current state-of-the-art image generation model can produce images with strange details, blurriness, and unrealistic physical phenomena such as lighting issues. These issues limit the model's ability to generate images with high accuracy and fidelity to real-world scenes. Overall, while the SOTA image generation model has shown remarkable capabilities, it still faces significant challenges that need to be addressed for it to achieve even greater success in the field of image generation.

## F.2 Broader Impact

**Societal risks.** As AIGC continues to be promoted in various fields, concerns about its societal use have become increasingly prominent. These concerns involve various issues such as bias and ethics. As we have demonstrated, it is getting more and more difficult for humans to distinguish between AI-generated images and real images. Therefore, AI models may produce content that contradicts or even absurdly violates reality, posing a serious threat to factual information. Photos may then become increasingly difficult to use as evidence in the future, and even serious public opinion effects may result. For example, there were many AI-generated images of Trump being arrested on Twitter recently [3]. Such content may be used to spread false information, incite violence, or harm individuals or organizations. Besides, AIGC can be used to create realistic virtual characters, which may be used for malicious purposes such as online fraud, scams, or harassment.

It is crucial for researchers and practitioners in the field of AIGC to develop strategies to mitigate potential negative impacts. This includes developing methods to identify AI-generated images, establishing guidelines for their ethical use, and raising public awareness about their existence and potential impact. Only by working together can we ensure that the benefits of AIGC are fully realized while minimizing its negative consequences.

**Positive impacts.** Given that AI has shown remarkable performance in creating works of art and photography, it is expected to have a significant impact on artists and photographers in the real world. People can obtain a large number of desired works or photos at a lower cost, which could compress the market for artists and photographers. In this era of fast-food images, where should the new generation of artists and photographers go [2]? However, AI can only generate soulless works, lacking the creativity, imagination, and emotion possessed by human artists and photographers [51]. Even the most advanced AI technology cannot replace the creativity and individuality of human artists and photographers. Therefore, although the emergence of AI has indeed brought new challenges and changes to the fields of art and photography, human artists and photographers are still highly valued.

The emergence of AI technology presents various new opportunities for artists, designers, and users. One of the most significant benefits is the ability to create new and innovative visual works, such as digital art and logos, while reducing the time and cost associated with traditional image creation methods. AI technology allows people to generate unique and novel images that might not have been possible otherwise, leading to new ideas and inspiration. Moreover, AI technology can help optimize existing works of art and photos, leading to improved quality and value. For instance, AI can be used to enhance or restore old or damaged images to their original state [77], which can be particularly useful in restoring historic photographs or artworks [9]. AIGC also provides users a more personalized experience by creating images tailored to their personal preferences [28, 65]. This customization can lead to more engaging and immersive experiences for users.

**Academic impacts.** In this study, we conduct a quantitative human evaluation of whether the most advanced AI model can deceive the human eye. Results indicate several academic directions that could be explored in the future:

• Since people cannot discern the authenticity of images, a natural question arises: *Can AI distinguish whether an image is generated by AI?* Exploring how to use AI to detect AI-generated images is a problem that could be studied [79]. Establishing a detection system to recognize AIGC will greatly ensure the security of society and the credibility of images.

• Even the most advanced image generation model still cannot guarantee the stable generation of high-quality images. At the same time, as shown in Fig. 9, AI-generated images often have certain defects. Our failure case analysis will inspire researchers to design better image generation models. Exploring how to solve these AIGC defects is an important future research direction.

• There is an interesting phenomenon in MPBench: CLIP-ViT-L (LC) [57] freezes the pre-trained backbone and unfreezes the last linear layer. Its generalization in MPBench is very good, but its accuracy in real images has dropped a lot. However, other models initialized from pre-trained models with whole backbone unfrozen have good accuracy in real images, but the generalization in MPBench are not good. This phenomenon shows an interesting research problem: Can we achieve a balance between these two settings? To study how many proportions of backbones should be frozen and how many proportions of backbones should be unfrozen is the best setting for fake image detection task is a good research problem.

• In the real world, it is difficult to obtain comprehensive and diverse data, leading to the famous problem of data imbalance [37]. Using imbalanced data will result in various issues such as the long-tail problem [85] and bias problem [53]. Since the current state-of-the-art image generation model can already produce high-quality data, exploring how to use the image generation model to solve these problems and test the current model's robustness and bias is a problem that could be studied.

## F.3  Limitation

While this work has so far provided several state-of-the-art and large-scale training and validation datasets, as well as several powerful benchmarks, this section explores the limitations of the which are expected to be addressed in future studies.

**Dataset limitation.**   Our training dataset Fake2M only includes three advanced models: Stable Diffusion v1.5 [64], IF [5], and StyleGAN3 [41], limited by the absence of open-source and powerful open vocabulary GAN [38] and Autoregressive models [81]. Due to the lack of API, we are unable to provide a training dataset for Midjourney V5. We hope that future work can further improve the diversity and size of the training dataset to include more powerful generative models.
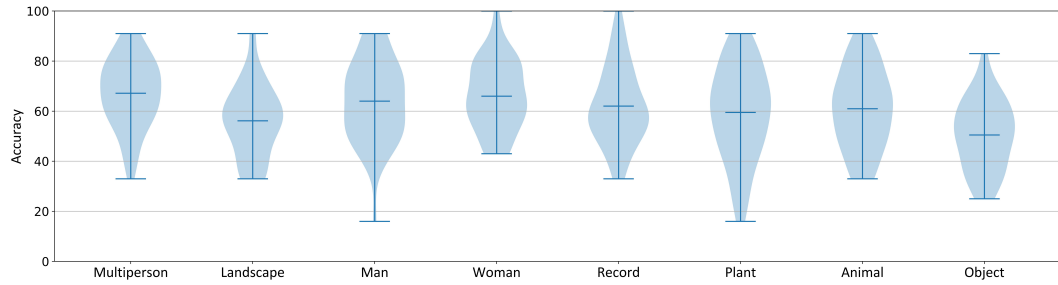
For the validation datasets, we only include validation datasets for the most advanced generative models, without including validation datasets for other tasks, such as deepfake and low-level tasks. We hope that future work can further improve the diversity of the validation dataset to include more tasks about fake images.

**Benchmark limitation.**   Due to the resource limitations, our high-quality human evaluation HP-Bench only recruits 50 participants. Our human evaluation also lacks diversity in terms of participant background, as it only includes a few attributes such as age, AIGC-background and gender. We hope that future work can further improve the diversity and size of the participants.
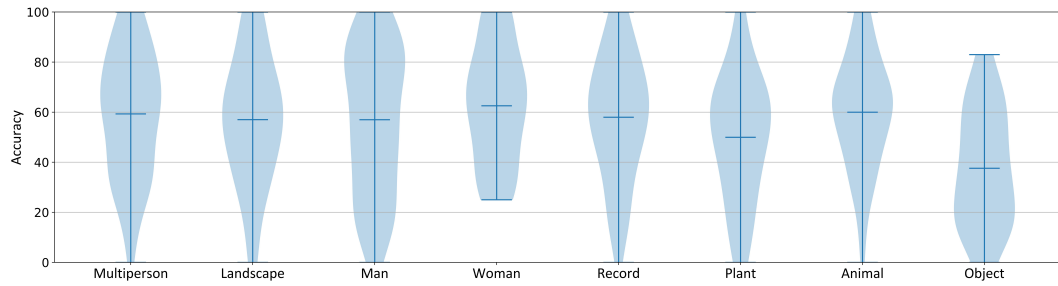
## F.4  Conclusion

In this study, we present a comprehensive evaluation of both human discernment and contemporary AI algorithms in detecting fake images. Our findings reveal that humans can be significantly deceived by current cutting-edge image generation models: high-quality AI-generated images can be comparable to real photographs. In contrast, AI fake image detection algorithms demonstrate a superior ability to distinguish authentic images from fakes. Despite this, our research highlights that existing AI algorithms, with a considerable misclassification rate of 13%, still face significant challenges. We anticipate that our proposed dataset, **Fake2M**, and our dual benchmarks, **HPBench** and **MPBench**, will invigorate further research in this area and assist researchers in crafting secure and reliable AI-generated content systems. As we advance in this technological era, it is crucial to prioritize responsible creation and application of generative AI to ensure its benefits are harnessed positively for society.
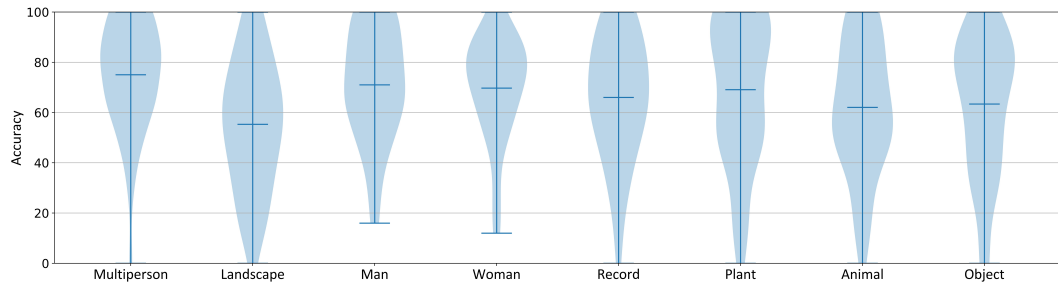
We have focused on the surprising abilities of the current SOTA image generation model, but we have not addressed the core questions of why and how it achieves such remarkable intelligence, nor the most important issues of how to ensure the security and credibility of AIGC images. It is a significant challenge for researchers to develop secure and reliable AIGC systems that can be trusted for various real-world applications, and ensure the responsible and ethical use of AIGC technology in the future. It is time to prioritize responsible development and the use of generative AI to ensure a positive impact on society.

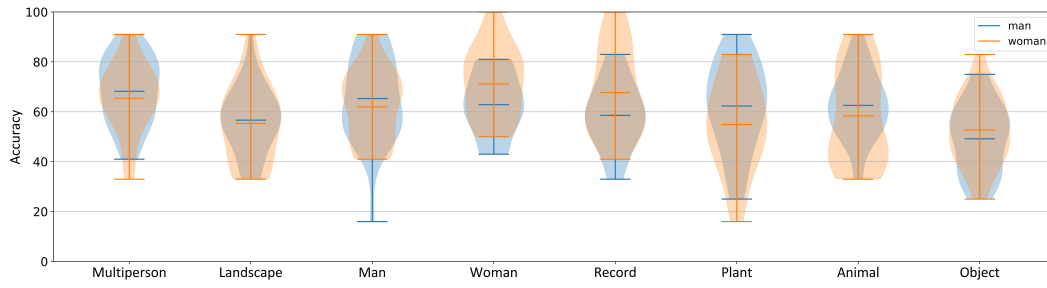(a) All images with different categories for all persons

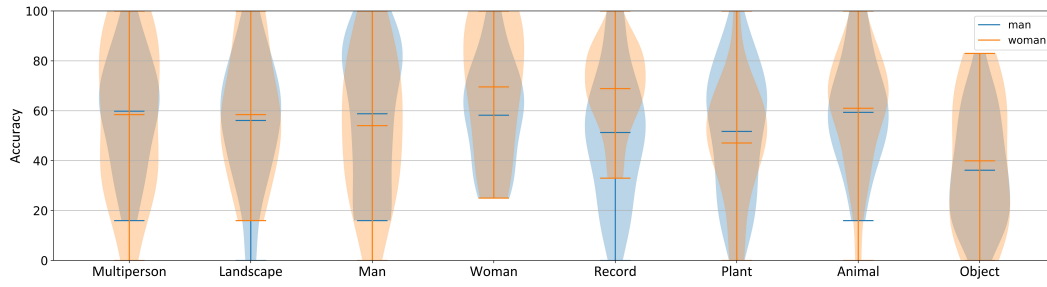(b) Only AI-generated images with different categories for all persons

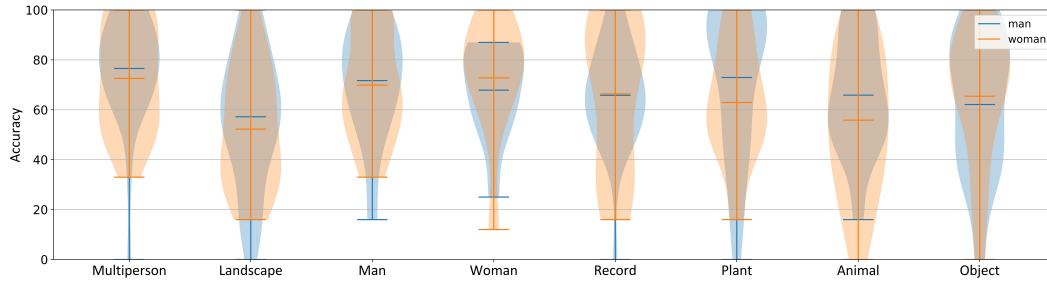(c) Only real images with different categories for all persons

Figure 12: **Score distributions for all volunteers with different categories.**

(a) All images with different categories for man and woman

(b) Only AI-generated images with different categories for man and woman

(c) Only real images with different categories for man and woman

Figure 13: **Score distributions for men and women with different categories.**
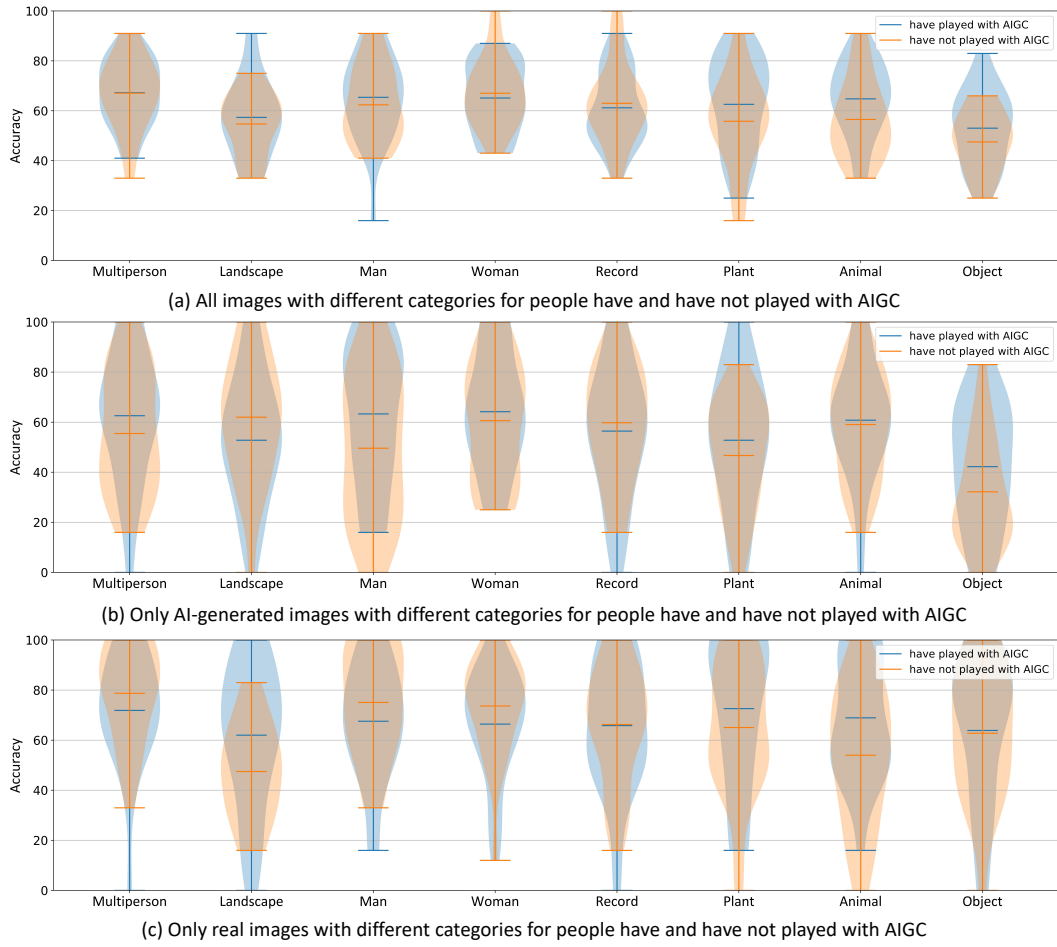
(a) All images with different categories for people have and have not played with AIGC

(b) Only AI-generated images with different categories for people have and have not played with AIGC

(c) Only real images with different categories for people have and have not played with AIGC

Figure 14: **Score distributions for volunteers with and without AIGC background.**