

The Visual Counter Turing Test (VCT²): A Benchmark for Evaluating AI-Generated Image Detection and the Visual AI Index (V_{AI})

Anonymous ACL submission

Abstract

The rapid progress and widespread availability of text-to-image (T2I) generation models have heightened concerns about the misuse of AI-generated visuals, particularly in the context of misinformation campaigns. Existing AI-generated image detection (AGID) methods often overfit to known generators and falter on outputs from newer or unseen models. To systematically address this generalization gap, we introduce the **Visual Counter Turing Test (VCT²)**, a comprehensive benchmark of 166,000 images, comprising both real and synthetic prompt-image pairs produced by six state-of-the-art (SoTA) T2I systems: Stable Diffusion 2.1, SDXL, SD3 Medium, SD3.5 Large, DALL·E 3, and Midjourney 6. We curate two distinct subsets: *COCO_{AI}*, featuring structured captions from MS COCO, and *Twitter_{AI}*¹, containing narrative-style tweets from The New York Times. Under a unified zero-shot evaluation, we benchmark 17 leading AGID models and observe alarmingly low detection accuracy, 58% on *COCO_{AI}* and 58.34% on *Twitter_{AI}*. To transcend binary classification, we propose the **Visual AI Index (V_{AI})**, an interpretable, prompt-agnostic realism metric based on twelve low-level visual features, enabling us to quantify and rank the perceptual quality of generated outputs with greater nuance. Correlation analysis reveals a moderate inverse relationship between V_{AI} and detection accuracy: Pearson ρ of -0.532 on *COCO_{AI}* and ρ of -0.503 on *Twitter_{AI}*; suggesting that more visually realistic images tend to be harder to detect, a trend observed consistently across generators. We release *COCO_{AI}* and *Twitter_{AI}* to catalyze future advances in robust AGID and perceptual realism assessment.

1 Introduction

The rapid advancement of text-to-image (T2I) generative models, such as Stable Diffusion (Rom-

¹Midjourney 6 is excluded from *Twitter_{AI}* due to prompt filtering constraints.



Figure 1: An AI-generated image of Pope Francis wearing a gigantic white puffer jacket went viral on social media platforms like Reddit and Twitter (X) in March 2023. This image sparked widespread media discussions on the potential misuse of generative AI technologies, becoming an iconic example of AI-generated misinformation. For more details, see the [Forbes story](#).

bach et al., 2022; Podell et al., 2023; Esser et al., 2024), DALL·E (Ramesh et al., 2021, 2022; Betker et al., 2023), Midjourney (Midjourney, 2024), and Imagen (Saharia et al., 2022), has revolutionized visual content creation. These models unlock powerful creative workflows and democratize image synthesis at scale. However, their widespread accessibility also raises critical concerns about visual misinformation and content authenticity. As illustrated in Figure 1, synthetic images can convincingly mimic journalistic or photographic style, blurring the boundary between real and generated content. This growing threat has prompted global attention. In March 2023, an open letter (Marcus, 2023) warned that generative AI could destabilize the global information ecosystem. The European Commission reported a significant decline in online content moderation accuracy, from 90.4% in 2020 to just 64.4% in 2022 (Commission, 2022). Meanwhile, social platforms process over 3.2 billion images and 720,000 hours of video daily (T.J. Thomson, 2020), with synthetic media projected to account for 90% of online content by 2026 (Europol,

2024).

Despite increasing demand for reliable detection tools, existing AI-generated image detection (AGID) methods often fail to generalize to images from unseen generators or real-world contexts. Watermark-based approaches remain fragile, easily circumvented via cropping, filtering, or adversarial manipulation (Zhao et al., 2025). Meanwhile, prior AGID benchmarks (Zhu et al., 2023; Sha et al., 2023) suffer from limited real-image diversity, narrow prompt coverage, outdated model inclusion, and closed access, impeding rigorous evaluation and progress.

To address these limitations, we introduce the **Visual Counter Turing Test (VCT²)**, a large-scale benchmark dataset for zero-shot AGID evaluation. VCT² contains approximately 166,000 images, including 26,000 real prompt-image pairs and 140,000 synthetic images produced by six SoTA T2I models: Stable Diffusion 2.1, SDXL, SD3 Medium, SD3.5 Large, DALL-E 3, and Midjourney 6, spanning both open-source and proprietary systems. The prompts in VCT² are drawn from two semantically distinct sources to capture both structured and open-ended language. The COCO_{AI} subset uses object-centric captions from MS COCO (Lin et al., 2014), a staple in vision-language research. The Twitter_{AI} subset comprises narrative-style tweets authored by The New York Times (@nytimes), providing real-world, journalistic prompts rich in nuance and context. This diversity allows us to evaluate AGID methods across a wide range of generation styles and domains.

To enable more nuanced evaluation beyond binary classification, we introduce the **Visual AI Index (V_{AI})**, a model-agnostic, interpretable metric that quantifies the perceptual realism of an image based solely on its visual content. V_{AI} produces a scalar score derived from twelve handcrafted, low-level image features, including texture complexity, frequency-domain statistics, Haralick features, and image sharpness. These features have been selected based on their empirically observed alignment with human judgments of realism. Correlation analysis further supports the utility of V_{AI} as a proxy for detection difficulty: we observe a moderate inverse relationship between V_{AI} scores and AGID detection accuracy across generative models (Pearson $\rho = -0.503$ on Twitter_{AI} and $\rho = -0.532$ on COCO_{AI}), indicating that more visually realistic images tend to be harder to detect. Our realism scores offer a prompt and model-agnostic lens into

the perceptual quality of generated images.

We evaluate **17 AGID** methods under a standardized zero-shot setting, using publicly available implementations and default model checkpoints. Our goal is to assess how well these methods generalize across a variety of text-to-image models, including open-source systems like Stable Diffusion 2.1, SDXL, SD3 Medium, and SD3.5 Large, as well as proprietary models such as DALL-E 3 and Midjourney 6, and across two domains: structured and high quality MS COCO captions and images and narrative-style tweets and real world images from The New York Times. Experimental results (Section 4) reveal model generalization gaps by noticeable detection performance degradation, with average detection accuracy of 58% on COCO_{AI} and 58.34% on Twitter_{AI}. We observe lower detection accuracy on COCO_{AI} compared to Twitter_{AI}. This is likely because COCO prompts produce images that are more photo-realistic and visually similar to real photos. In contrast, Twitter_{AI} generations often include creative or unusual visual patterns, leading to more detectable differences. Notably, DALL-E 3 and SD3.5 consistently yield the lowest detection accuracy across both domains. To summarize, our main contributions are:

(i) We introduce the **Visual Counter Turing Test (VCT²)** benchmark to evaluate the generalization capabilities of AI-generated image detection methods across diverse prompt styles and real image sources, including MS COCO and Twitter, as well as six state-of-the-art synthetic image generators.

(ii) We define the **VisualAI Index (V_{AI})**, a scalar metric to quantify perceptual realism based on twelve interpretable low-level visual features.

2 Recent Advances in AI-Generated Image Detection Techniques

AI-generated image detection (AGID) is becoming increasingly vital as synthetic content continues to grow in both photorealism and scale. Detection methods vary widely in their design assumptions, feature representations, and robustness to distribution shifts, such as changes in generative models, prompt styles, or real image characteristics that diverge from photographic norms.

To facilitate systematic evaluation, we categorize AGID approaches into three broad groups:

(i) **Generation Artifact-Based Methods:** These methods target low-level signals introduced

during the image synthesis process, such as upsampling artifacts, denoising residuals, or color inconsistencies. While often computationally efficient, they tend to be fragile under post-processing or model variation.

(ii) **Feature Representation-Based Methods:**

These rely on high-level semantic or perceptual features extracted via CNNs, Vision Transformers, or CLIP-style encoders. They typically offer stronger generalization across domains, though may miss fine-grained generative artifacts.

(iii) **Hybrid Methods:** These approaches integrate both low- and high-level cues, and often leverage contrastive learning, multi-modal embeddings, or text-image alignment to enhance robustness under distributional shifts.

Figure 2 illustrates this taxonomy. We evaluate 17 publicly available AGID models spanning all three categories, selected based on their methodological diversity, recent relevance, and open-source availability. Each method is tested in a standardized zero-shot setting using its default checkpoint, without any fine-tuning on the VCT² benchmark. This taxonomy provides a reference framework for interpreting detection trends discussed in Section 4, with further implementation details outlined in Appendix B.

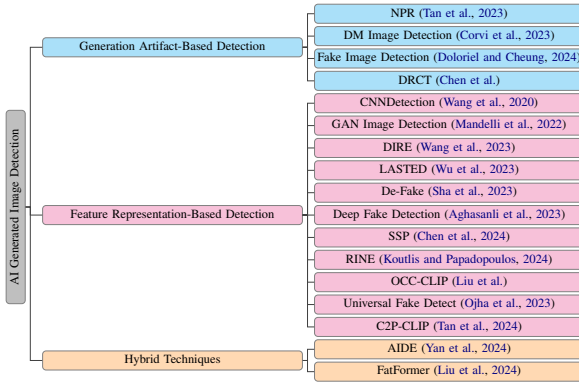


Figure 2: The taxonomy of AI-generated image detection techniques, categorized into three main groups: Generation Artifact-Based Detection, Feature Representation-Based Detection, and Hybrid Techniques.

3 The Visual Counter Turing Test (VCT²) Benchmark Dataset

We introduce the Visual Counter Turing Test (VCT²), a large-scale benchmark designed to eval-

uate AI-generated image detection (AGID) techniques. VCT² includes the following:

- Approximately 26,000 real image-prompt pairs, combining our curated Twitter dataset and the benchmark MS COCO dataset;
- Approximately 140,000 synthetic images, generated using six state-of-the-art text-to-image models; open-source models: Stable Diffusion 2.1, SDXL, SD3 Medium, SD3.5 Large; and two proprietary models: DALL-E 3, Midjourney 6.
- In total, around 166,000 images derived from 26,000 unique prompts.

This scale provides a balance between structured, caption-based content and naturalistic, real-world prompts, positioning VCT² among the most comprehensive AGID datasets to date.

3.1 Prompt Sources and Coverage

To ensure diversity in both semantic content and visual generation styles, we curated prompts from two distinct and complementary sources:

- ~10,000 benchmark prompts from the MS COCO dataset (Lin et al., 2014), focused on object-centric and everyday scenes;
- ~16,000 real-world prompts from the @nytimes Twitter account 2011–2023. To assess topical diversity, we identified ten topics and associated keywords, and then assigned each tweet to its most probable topic. Table 1 shows the five dominant clusters. These clusters reflect both editorial depth and real-world content breadth. Their presence enhances the semantic realism of our benchmark and supports rigorous AGID evaluation across multiple domains.

3.2 Real Twitter Prompt-Image Dataset Collection and Processing

To construct a diverse and reliable dataset of real Twitter images, we employed an automated data collection pipeline using Python and Selenium. We focused on tweets from @nytimes (The New York Times) due to its editorial credibility, rigorous fact-checking, and diverse topical coverage.

Data Collection. Our pipeline sampled tweets spanning a 12-year period (2011–2023), retaining

Table 1: Topic Clusters in the NYT Twitter Subset.

Topic Cluster	Tweet Count	Top Keywords
Daily Briefings and News Summaries	1129	<i>know, need, day, morning, briefing, evening</i>
New York City and Culture	1961	<i>new, york, city, times, books, critics</i>
Art, Movies, and Obituaries	631	<i>photo, review, obituary, art, movie, critic</i>
Health, COVID-19, and Breaking News	1436	<i>coronavirus, health, opinion, news, breaking, people</i>
Opinion Pieces and Societal Reflections	914	<i>nytopinion, life, young, america, death, ebola</i>
Travel and International Destinations	605	<i>hours, italy, florida, china, japan, park</i>
World Events and Sports	903	<i>world, cup, photos, team, war, country</i>
Lifestyle and City Aesthetics	1164	<i>like, looks, look, city, love, idea</i>
Time, Life Stories, and Incarceration	1200	<i>years, life, ago, prison, time, close</i>
Home, Food, and Leisure	965	<i>make, recipes, summer, home, simple, best</i>
miscellaneous	5001	–

only those with attached media. The goal was to align real images with captions that could feasibly be used to generate synthetic counterparts.

Definition of Real Images. We define “real” images as those not generated by AI. This includes natural photographs as well as editorial media such as UI screenshots, infographics, and photojournalistic illustrations, provided they are not produced using generative models. This definition reflects the ambiguity present in real-world detection scenarios, where non-photographic content may still be authentic.

Data Filtration and Preprocessing. To ensure quality and consistency, we applied several filtering steps: (i) removal of duplicate tweets and media; (ii) exclusion of irrelevant content such as word games or puzzles; and (iii) filtering of non-English tweets. Additionally, preprocessing involved removal of hashtags and URLs, and retention of only alphanumeric characters to facilitate downstream analysis and clustering. All real images and prompts are organized into two structured subsets: *COCO_{AI}* and *Twitter_{AI}*, which are publicly released.

3.3 Benchmark Scale and Contributions

VCT² offers several key advantages over existing benchmarks:

- **Scale.** VCT² contains 166,000 images generated from 26,000 unique prompts. While GenDet (Zhu et al., 2023) includes a larger total number of images (770,000), it only contains 70,000 real images and is less balanced across prompt domains. De-Fake (Sha et al., 2023), in contrast, includes approximately 100,000 images with more limited generator and prompt coverage.
- **Model diversity.** GenDet and De-Fake focus primarily on earlier open-source models such

as BigGAN, StyleGAN2, and Stable Diffusion v1.4. VCT² includes six cutting-edge generative models supporting broader evaluation across current-generation image generators.

- **Prompt realism.** VCT² uniquely combines benchmark-style prompts from MS COCO with naturalistic, real-world prompts curated from a 12-year archive of @nytimes tweets, capturing diverse linguistic styles and topics.
- **Mixed-media realism.** The real image subset includes ambiguous formats such as infographics, UI screenshots, and editorial photos, reflecting the heterogeneous content encountered in real-world detection scenarios.
- **Public accessibility.** All prompts, real and synthetic images, and evaluation scripts for 17 AGID baselines are publicly released to facilitate reproducibility and comparative benchmarking.

To our knowledge, VCT² is the first large-scale AGID benchmark to pair real-world journalistic prompts with diverse state-of-the-art text-to-image models, providing a robust and publicly available testbed for evaluating both detection performance and perceptual realism across different prompt domains and generative model types.

4 Evaluation and Results

We evaluate the VCT² benchmark under zero-shot settings using 17 state-of-the-art AI-generated image detection (AGID) methods. These span artifact-, feature-, and hybrid-based approaches. In the following, we present detection performance, examine cross-domain and cross-model generalization, and analyze detector sensitivity to detector type.

4.1 Evaluation Protocol

To simulate real-world deployment, we assess all detectors without fine-tuning. Public checkpoints and default hyperparameters are used. Performance is measured separately on COCO_{AI} and Twitter_{AI} subsets, reporting accuracy, precision, and recall. Results are summarized in Tables 2 and 3.

4.2 Cross-Domain and -Model Trends

Figure 3 presents the average detection accuracies per generator across the two domains. Overall, detection performance is low, with accuracy dropping further on COCO_{AI} compared to Twitter_{AI}.

Detection performance also varies across generators. Images from earlier models like SD2.1 and SDXL remain relatively detectable. In contrast, newer or proprietary models such as SD3.5 Large and DALL-E 3 yield significantly lower detection results, suggesting that existing detectors may be overfitted to older, synthetic image distributions.

4.3 Comparative Detector Performance

Feature-based detectors (e.g., De-Fake (Sha et al., 2023)) and hybrid methods (e.g., DRCT (Chen et al.)) generally outperform artifact-based detectors like CNNDetection (Wang et al., 2020) and NPR (Tan et al., 2023). The latter collapse on proprietary models due to reliance on low-level artifacts often absent in advanced generators.

Conversely, feature-based and contrastive methods benefit from semantic representations, allowing stronger generalization to unseen prompt styles and model outputs. DRCT (ConvB and UnivFD) and De-Fake show consistent robustness across both subsets.

5 The Visual AI Index (V_{AI})

We introduce the *Visual AI Index* (V_{AI}), an interpretable, prompt-agnostic metric that scores the perceptual realism of images based on low-level visual features. V_{AI} provides a continuous score that reflects where an image lies along a spectrum

Table 2: Overall accuracy (Acc), recall (R), and precision (P) across COCO_{AI} synthetic datasets generated from MS COCO prompts. All values are in %. Color-coded: Green ($\geq 90\%$), Yellow-Green (80–89%), Yellow (70–79%), Orange (60–69%), red ($<60\%$).

Method	SD2.1			SDXL			SD3 Medium			SD3.5 Large			DALL-E 3			Midjourney 6		
	Acc	R	P	Acc	R	P	Acc	R	P	Acc	R	P	Acc	R	P	Acc	R	P
CNNDetection (Wang et al., 2020)	49.94	0.03	65.11	49.96	0.07	77.52	49.93	0.01	81.16	49.99	0.14	33.04	49.93	0.00	35.13	49.95	0.05	63.15
NPR (Tan et al., 2023)	26.76	1.89	34.26	26.68	1.73	33.15	27.96	4.29	34.41	70.32	48.37	79.44	25.81	0.00	41.13	25.81	0.00	48.13
DM Image Detection (Corvi et al., 2023)	83.92	67.92	99.40	69.96	40.00	98.91	63.58	27.23	98.04	38.58	32.06	0.07	49.96	0.00	40.00	51.73	3.52	87.04
Fake Image Detection (Doloriel and Cheung, 2024)	49.84	0.49	63.58	49.83	0.48	66.68	50.02	0.86	66.91	48.24	0.11	62.57	49.59	0.00	34.90	49.79	0.40	62.89
DIRE (Wang et al., 2023)	47.08	93.40	37.66	49.67	98.57	47.07	48.59	96.40	38.88	50.63	99.23	58.68	48.89	97.01	43.25	50.04	99.31	52.74
LASTED (Wu et al., 2023)	54.00	8.67	56.62	61.13	9.86	61.20	51.87	9.61	57.67	55.21	10.11	57.35	66.18	44.85	76.21	68.21	14.37	63.14
GAN Image Detection (Mandelli et al., 2022)	51.87	82.93	51.16	56.35	91.75	53.72	58.26	95.35	54.74	45.61	79.08	47.37	48.10	74.93	48.77	57.15	93.42	54.14
AIDE (Yan et al., 2024)	60.30	20.98	93.77	64.34	28.91	96.75	57.11	14.45	94.28	50.83	5.01	52.31	50.00	0.02	61.23	76.01	52.25	96.92
SSP (Chen et al., 2024)	50.15	99.63	50.07	49.95	99.63	49.97	50.34	99.63	50.17	50.30	99.48	50.29	49.91	99.63	49.95	49.95	99.63	49.97
FatFormer (Liu et al., 2024)	50.00	0.00	0.00	50.00	0.00	0.00	50.00	0.01	100	50.28	0.00	0.00	48.01	0.00	0.00	48.01	0.00	0.00
DRCT (ConvB) (Chen et al.)	98.76	99.61	97.94	96.83	95.75	97.86	80.72	63.54	96.81	78.58	59.05	96.51	49.99	2.08	49.76	67.48	37.06	94.65
DRCT (UnivFD) (Chen et al.)	88.57	96.98	83.02	89.45	98.73	83.27	84.90	89.64	81.88	83.09	84.09	82.29	79.98	79.80	80.09	89.64	99.12	83.32
RINE (Koutlis and Papadopoulos, 2024)	74.43	49.63	98.49	56.47	13.71	94.76	61.99	24.75	97.03	55.34	27.72	95.34	50.05	0.87	53.37	63.13	27.02	97.27
OCC-CLIP (Liu et al.)	51.49	92.28	50.82	47.11	14.95	41.91	50.60	66.03	50.46	49.08	50.67	67.63	78.82	50.28	55.04	75.04	53.60	50.03
De-Fake (Sha et al., 2023)	92.37	97.90	88.15	91.23	95.62	87.90	91.30	95.76	87.92	52.57	86.05	5.11	90.58	94.31	87.76	86.22	85.59	86.68
Deep Fake Detection (Aghasanli et al., 2023)	49.49	49.49	49.03	51.43	51.43	49.65	49.85	49.85	49.97	50.66	50.66	52.19	52.73	52.73	53.02	52.87	52.87	54.09
Universal Fake Detect (Ojha et al., 2023)	74.42	77.15	73.15	69.18	65.84	70.56	70.11	68.79	70.66	57.46	60.10	57.09	50.00	99.99	50.00	53.23	76.28	52.21
C2P-CLIP (Tan et al., 2024)	53.38	7.53	90.73	53.69	8.15	91.30	55.50	11.77	93.87	52.13	4.40	97.01	49.76	0.29	27.36	50.31	1.40	64.52

Table 3: Overall accuracy (Acc), recall (R), and precision (P) across Twitter_{AI} synthetic datasets generated from Twitter prompts. Midjourney 6 is not included as it blocks image generation for most Twitter prompts. All values are in %. Color-coded: Green ($\geq 90\%$), Yellow-Green (80–89%), Yellow (70–79%), Orange (60–69%), red ($<60\%$).

Method	SD2.1			SDXL			SD3 Medium			SD3.5 Large			DALL-E 3		
	Acc	R	P	Acc	R	P	Acc	R	P	Acc	R	P	Acc	R	P
CNNDetection (Wang et al., 2020)	50.00	0.06	52.21	49.98	0.03	59.98	50.19	0.44	74.35	50.34	0.76	76.04	49.97	0.01	34.59
NPR (Tan et al., 2023)	50.23	2.22	50.89	50.46	2.68	60.58	51.45	4.66	68.26	52.12	7.81	67.44	49.12	0.00	42.20
DM Image Detection (Corvi et al., 2023)	88.31	77.57	97.82	73.82	48.58	93.74	65.15	31.24	90.34	63.56	28.19	89.66	49.53	0.00	33.34
Fake Image Detection (Doloriel and Cheung, 2024)	49.86	0.53	56.33	49.88	0.58	60.83	50.35	1.51	66.15	48.57	1.12	62.13	49.59	0.01	33.11
DIRE (Wang et al., 2023)	43.90	86.95	36.20	48.57	96.29	46.29	48.49	96.13	38.48	49.41	98.01	45.11	46.33	91.81	36.19
LASTED (Wu et al., 2023)	77.60	1.93	59.60	83.60	2.75	66.04	83.24	2.75	61.52	82.71	2.59	61.90	78.77	25.57	76.81
GAN Image Detection (Mandelli et al., 2022)	53.26	77.37	52.25	55.84	82.36	53.86	60.01	91.04	56.21	54.78	80.60	53.19	53.99	79.44	52.68
AIDE (Yan et al., 2024)	55.69	11.81	81.98	60.43	21.29	89.61	56.49	13.41	87.40	56.57	6.16	58.42	49.93	0.25	43.61
SSP (Chen et al., 2024)	49.91	99.66	49.95	50.20	99.66	50.10	50.20	99.66	50.10	54.94	99.33	55.04	50.18	99.66	50.10
FatFormer (Liu et al., 2024)	50.04	0.08	100	50.04	0.08	100	50.00	0.00	100	55.10	0.02	100	50.02	0.00	0.00
DRCT (ConvB) (Chen et al.)	96.81	99.77	94.20	93.96	94.05	93.87	71.79	49.73	89.01	77.87	59.54	90.35	47.31	0.76	11.02
DRCT (UnivFD) (Chen et al.)	67.47	96.73	61.02	68.32	98.43	61.44	64.81	91.40	59.67	62.94	85.60	55.68	53.76	69.30	52.87
RINE (Koutlis and Papadopoulos, 2024)	77.07	55.40	97.79	57.86	16.97	93.13	62.13	25.50	95.32	66.36	44.37	94.36	49.61	0.48	27.64
OCC-CLIP (Liu et al.)	46.88	74.11	47.98	45.67	51.17	46.10	48.84	67.54	49.16	47.82	66.16	48.03	47.75	45.63	49.72
De-Fake (Sha et al., 2023)	81.13	91.51	75.78	78.16	85.57	74.53	79.39	88.03	72.06	40.80	0.00	0.00	79.95	89.14	75.29
Deep Fake Detection (Aghasanli et al., 2023)	50.80	50.80	51.84	53.64	53.64	56.59	51.44	51.44	51.51	49.19	49.19	56.38	55.30	55.30	60.34
Universal Fake Detect (Ojha et al., 2023)	72.88	74.17	72.31	69.47	73.91	67.89	68.41	72.58	67.00	55.38	45.60	56.69	50.00	99.99	50.00
C2P-CLIP (Tan et al., 2024)	52.21	6.74	88.76	53.28	7.98	91.25	47.92	0.97	26.17	54.16	8.47	98.39	49.73	0.21	53.45

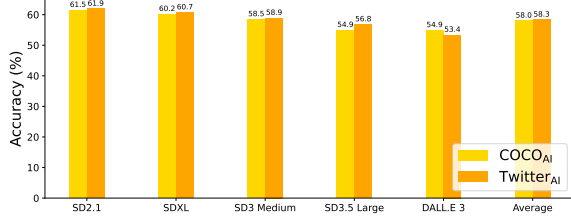


Figure 3: Average detection accuracy across the COCO_{AI} and Twitter_{AI} subsets for each generator.

of visual realism. Many real images in web-scale datasets (e.g., news media, social platforms) are not pristine photographs; they may include screenshots, digital graphics, or compressed visuals. These images often lack sharpness, contrast, or structure. V_{AI} quantifies perceptual quality by learning to score realism using a combination of handcrafted visual features, independent of prompts or model-specific information.

5.1 Feature Design

V_{AI} uses twelve visual features grouped into three categories:

(i) **Texture and Frequency:** Texture Complexity, Haralick Contrast, Haralick Correlation, Haralick Energy, Frequency Mean, Frequency Standard Deviation.

(ii) **Sharpness and Structure:** Image Sharpness, Image Smoothness, Image Contrast.

(iii) **Color and Semantics:** Color Distribution Consistency, Object Coherence, Contextual Relevance.

Texture Complexity quantifies the variety and unpredictability of an image’s texture. It is determined by computing the entropy of the normalized Local Binary Pattern (LBP) histogram of the grayscale image using the formula $-\sum_{k=0}^{P-1} \tilde{H}_{LBP}(k) \log_2(\tilde{H}_{LBP}(k) + \epsilon)$. Here, $\tilde{H}_{LBP}(k)$ represents the normalized histogram value for LBP bin k , and P is the total number of bins in the LBP histogram. The small constant ϵ (in our case, 1×10^{-6}) is used to avoid taking the logarithm of zero.

Haralick features are texture descriptors computed from the gray-level co-occurrence matrix (GLCM), which encodes the frequency $G(i, j)$ of pixel intensity pairs (i, j) occurring at a fixed spatial offset. We use three common features:

Haralick Contrast is defined as $\sum_{i,j} (i - j)^2 G(i, j)$, capturing local intensity variation.

Haralick Correlation is computed as $\sum_{i,j} \frac{(i - \mu_i)(j - \mu_j)G(i,j)}{\sigma_i \sigma_j}$, where μ_i, μ_j and σ_i, σ_j

are the means and standard deviations of the marginal GLCM distributions. It measures linear dependency between pixel pairs.

Haralick Energy (Angular Second Moment) is given by $\sum_{i,j} G(i, j)^2$, reflecting texture uniformity—higher values imply more homogeneous regions.

These values are averaged across multiple angles (e.g., $0^\circ, 45^\circ, 90^\circ, 135^\circ$) to ensure rotation-invariant descriptors.

We extract frequency-domain features using the 2D Fast Fourier Transform (FFT) of the grayscale image I . Let $\hat{I}(u, v) = \text{FFT2}(I)$ denote the Fourier-transformed image, and let $M(u, v) = |\hat{I}(u, v)|$ be the magnitude spectrum.

Frequency Mean is defined as $\text{FreqMean} = \frac{1}{HW} \sum_{u=1}^H \sum_{v=1}^W M(u, v)$, where $H \times W$ is the image resolution.

Frequency Standard Deviation is given by $\text{FreqStd} = \sqrt{\frac{1}{HW} \sum_{u=1}^H \sum_{v=1}^W (M(u, v) - \text{FreqMean})^2}$.

These two features capture the spectral energy and its variation. Higher values indicate detailed or noisy content, while lower values reflect smoother textures.

Image Sharpness is quantified as $\max(|I - I_{\text{blurred}}|)$. I and I_{blurred} denote the grayscale and blurred image with Gaussian kernel, respectively.

Image Smoothness evaluates how consistent the image’s texture is. It is quantified as $\frac{1}{1 + \text{var}(\Delta I)}$, where ΔI denotes the Laplacian of the grayscale image I .

Image Contrast measures the degree of variation in intensity across an image. It is quantified by calculating the standard deviation of the pixel values in the grayscale image, expressed as $\text{std}(I)$.

Color Distribution Consistency evaluates the variability in an image’s color distribution by analyzing the standard deviation of the normalized color histogram in the HSV color space. It is calculated as $\text{std}(\tilde{H}_{HSV}(h, s, v))$, where $\text{std}(\cdot)$ denotes the standard deviation of the normalized histogram $\tilde{H}_{HSV}(h, s, v)$ for hue h , saturation s , and value v .

Object Coherence evaluates the extent and clarity of edge detection in an image, providing insight into the consistency of object boundaries. It is determined using $\frac{\sum_{i,j} E(i,j)}{\sum_{i,j} 1}$, where $E(i, j)$ represents the value of the Canny edge image at pixel (i, j) , and the $\sum_{i,j} 1$ represents the total number of pixels in the image.

Contextual Relevance evaluates the distribution

of edge strengths across the image. It is given by $\text{var}(\sqrt{G_x^2 + G_y^2})$, where $\text{var}(\cdot)$ denotes the variance, and G_x and G_y are the gradients computed using the Sobel filter in the horizontal and vertical directions, respectively.

After Z-score normalization, each feature is standardized as $f_i(x) = \frac{v_i(x) - \mu_i}{\sigma_i}$, where $v_i(x)$ is the raw value of feature i for image x , and μ_i, σ_i are the mean and standard deviation of that feature across the dataset. The V_{AI} score is then computed as $V_{\text{AI}}(x) = \frac{s(x) - \min(s)}{\max(s) - \min(s)}$, where $s(x) = \sum_{i=1}^{12} w_i \cdot f_i(x)$ is the weighted sum of normalized features, and the min-max normalization rescales all scores $s(x)$ to the $[0, 1]$ range over the dataset.

To compute the Visual AI Index (V_{AI}), we learn a set of weights that quantify the contribution of each low-level feature to image realism. We use logistic regression to distinguish between real (label $y = 1$) and synthetic (label $y = 0$) images. Given a 12-dimensional feature vector $x = [f_1, f_2, \dots, f_{12}]$, the model estimates the probability that an image is real using the sigmoid function: $P(y = 1 | x) = 1/(1 + e^{-w^\top x})$, where w is the learned weight vector. The final Visual AI Index is defined as:

$$V_{\text{AI}}(x) = \frac{1}{1 + e^{-w^* \top x}},$$

where w^* is the optimized weight vector obtained after training. The weights are obtained by minimizing the binary cross-entropy loss: $\mathcal{L}(w) = \frac{1}{N} \sum_{i=1}^N [-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)]$, where $\hat{y}_i = 1/(1 + e^{-w^\top x_i})$. After training, we use the optimized weights w^* to compute the V_{AI} score. We train two models separately, one for COCO_{AI} and one for $\text{Twitter}_{\text{AI}}$, each tailored to the distribution of real images in its respective domain. Table 4 shows the final weights.

5.2 V_{AI} Analysis

We report the average V_{AI} scores for real and generated images across the COCO_{AI} and $\text{Twitter}_{\text{AI}}$ subsets in Figures 4 and 5. As expected, real images achieve the highest V_{AI} scores in both domains, reflecting the benchmark’s ability to assign higher realism to naturally occurring images. Among generative models, DALL-E 3 obtains the highest V_{AI} in both subsets (0.626 for COCO_{AI} , 0.593 for $\text{Twitter}_{\text{AI}}$), indicating its outputs most closely align with real images in terms of low-level features such as texture complexity, edge coherence, and color consistency. A cluster of diffusion-

Table 4: Learned V_{AI} feature weights for COCO_{AI} and $\text{Twitter}_{\text{AI}}$ domains.

Feature	COCO_{AI}	$\text{Twitter}_{\text{AI}}$
Texture Complexity	4.13	1.15
Color Dist. Consistency	-0.15	-0.05
Object Coherence	-0.87	1.56
Contextual Relevance	-1.33	2.42
Haralick Contrast	1.02	-4.52
Haralick Correlation	-0.42	-0.07
Haralick Energy	-1.46	-1.59
Freq. Mean	-0.87	-1.36
Freq. Std	-1.30	-0.20
Image Smoothness	0.37	-0.06
Image Sharpness	2.08	2.26
Image Contrast	0.50	-0.18

based models, SD2.1, SD3 Medium, and SD3.5 Large, follow DALL-E 3 with relatively similar V_{AI} scores, suggesting comparable levels of photo-realism. SDXL ranks lower across both domains, i.e. 0.496 COCO_{AI} and 0.573 $\text{Twitter}_{\text{AI}}$, which may be attributed to its tendency toward stylistic exaggeration or generation artifacts that deviate from natural image statistics. These artifacts can influence frequency-domain, edge-based, or texture descriptors negatively, despite the model’s high perceptual fidelity. Midjourney yields the lowest V_{AI} in the COCO_{AI} subset (0.432) and is excluded from the $\text{Twitter}_{\text{AI}}$ analysis due to the unavailability of corresponding generated images. We also report additional texture and semantic analyses in the Appendix: Local Binary Pattern (LBP) textures F.2, and PCA-based pairwise comparisons F.3 supporting our Visual AI Index (V_{AI}) findings on coherence and texture consistency across models.

5.3 Correlation with Detection Accuracy

To assess whether the Visual AI Index aligns with detection difficulty, we compute the Pearson correlation coefficient between the average V_{AI} scores and AGID detection accuracy across different generative models. The Pearson correlation is defined as $\rho = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$, where x_i and y_i are the V_{AI} score and AGID accuracy for model i , and \bar{x} and \bar{y} are their respective means. A higher value of ρ indicates stronger alignment between visual realism and detection performance.

To evaluate whether the Visual AI Index (V_{AI}) aligns with the difficulty of detecting AI-generated images, we compute the Pearson correlation coefficient between average V_{AI} scores and AGID detection accuracy across five generative models. The Pearson correlation coefficient ρ is defined as:

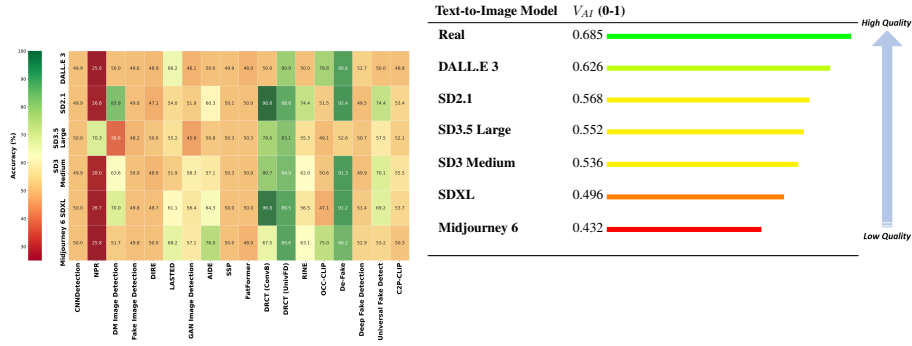


Figure 4: Right: V_{AI} scores of COCO_{AI} dataset. Left: Accuracy heat maps showing the average accuracy of each AGID method.

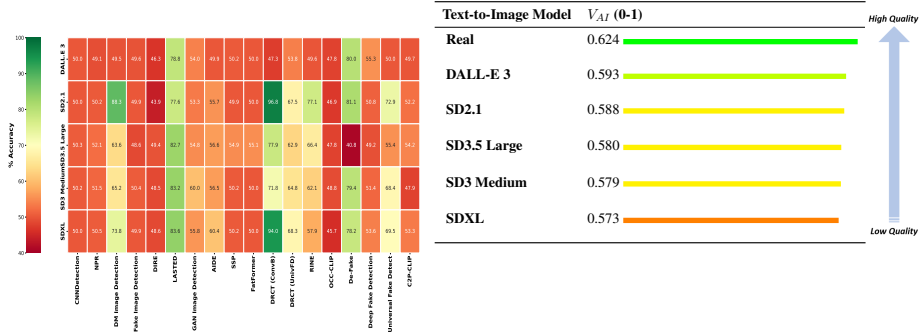


Figure 5: Right: V_{AI} scores of Twitter_{AI} dataset. Left: Accuracy heat maps showing the average accuracy of each AGID method.

$$\rho = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}},$$

where x_i and y_i represent the V_{AI} score and AGID accuracy for model i , and \bar{x} and \bar{y} are their respective means. The coefficient ρ ranges from -1 to 1 : a value near 1 implies a strong positive correlation, near -1 implies a strong negative correlation, and a value near 0 suggests no linear relationship. We compute ρ separately for the Twitter_{AI} and COCO_{AI} datasets. As shown in Table 5, our results indicate a moderate inverse relationship: models with higher visual realism tend to be harder to detect. However, the correlations are not statistically significant, likely due to the small number of generative models, i.e. $n = 5$, and should be interpreted cautiously.

6 Conclusion

In this paper, we introduced (VCT²), a comprehensive benchmark for evaluating AI-generated image

Table 5: Pearson correlation between V_{AI} and AGID detection accuracy.

Dataset	Pearson ρ	p-value
Twitter _{AI}	-0.503	0.388
COCO _{AI}	-0.532	0.356

detection (AGID) across diverse generative models, including cutting-edge proprietary systems like DALL-E 3 and Midjourney 6. By incorporating both real-world prompts and standardized captions, VCT² offers a challenging, realistic dataset for assessing generalization. The VCT² benchmark provides a critical resource for evaluating AGID techniques under challenging and varied conditions, highlighting performance gaps and guiding the development of more robust detection methods.

To assess the realism of images, we present the Visual AI Index (V_{AI}) that evaluates characteristics like texture complexity, Haralick correlation, frequency mean, and image sharpness. Our findings reveal that real images generally achieve higher V_{AI} scores than AI-generated images.

Limitations and Future Work

While our work provides a strong foundation for evaluating AGID methods and realism metrics, future directions include expanding to diverse domains (e.g., social platforms, synthetic video), integrating temporal and multimodal features into V_{AI} , and adapting it for localization or attribution. We also plan to explore human alignment and psychometric grounding of these continuous realism scores. As generative models evolve, updating the benchmark and exploring hybrid detection techniques will be key to ensuring resilience against increasingly sophisticated AI imagery.

References

- Agil Aghasanli, Dmitry Kangin, and Plamen Angelov. 2023. Interpretable-through-prototypes deepfake detection for diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 467–474.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, and 1 others. 2023. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8.
- Baoying Chen, Jishen Zeng, Jianquan Yang, and Rui Yang. Drct: Diffusion reconstruction contrastive training towards universal detection of diffusion generated images. In *Forty-first International Conference on Machine Learning*.
- Jiaxuan Chen, Jieteng Yao, and Li Niu. 2024. A single simple patch is all you need for ai-generated image detection. *arXiv preprint arXiv:2402.01123*.
- European Commission. 2022. [Eu code of conduct against online hate speech: latest evaluation shows slowdown in progress](#).
- Riccardo Corvi, Davide Cozzolino, Giada Zingarini, Giovanni Poggi, Koki Nagano, and Luisa Verdoliva. 2023. On the detection of synthetic images generated by diffusion models. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.
- Chandler Timm Doloriel and Ngai-Man Cheung. 2024. Frequency masking for universal deepfake detection. *arXiv preprint arXiv:2401.06506*.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, and 1 others. 2024. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*.
- Europol. 2024. [Facing reality: Law enforcement and the challenge of deepfakes](#). Accessed: 2024-08-30.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2021. Alias-free generative adversarial networks. *Advances in neural information processing systems*, 34:852–863.
- Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119.
- Christos Koutlis and Symeon Papadopoulos. 2024. Leveraging representations from intermediate encoder-blocks for synthetic image detection. *arXiv preprint arXiv:2402.19091*.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer.
- Fengyuan Liu, Haochen Luo, Yiming Li, Philip Torr, and Jindong Gu. Which model generated this image? a model-agnostic approach for origin attribution.
- Huan Liu, Zichang Tan, Chuangchuang Tan, Yunchao Wei, Jingdong Wang, and Yao Zhao. 2024. Forgery-aware adaptive transformer for generalizable synthetic image detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Sara Mandelli, Nicolò Bonettini, Paolo Bestagini, and Stefano Tubaro. 2022. Detecting gan-generated images by orthogonal training of multiple cnns. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 3091–3095. IEEE.
- Gary Marcus. 2023. [Pause giant ai experiments: An open letter](#).
- Midjourney. 2024. <https://www.midjourney.com/home>.
- Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. 2023. Towards universal fake image detectors that generalize across generative models. In *CVPR*.
- Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. 2023. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*.

675	Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey	Shilin Yan, Ouxiang Li, Jiayin Cai, Yanbin Hao, Xiao-	728
676	Chu, and Mark Chen. 2022. Hierarchical text-	long Jiang, Yao Hu, and Weidi Xie. 2024. A san-	729
677	conditional image generation with clip latents. <i>arXiv</i>	ity check for ai-generated image detection. <i>arXiv</i>	730
678	<i>preprint arXiv:2204.06125</i> , 1(2):3.	<i>preprint arXiv:2406.19435</i> .	731
679	Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott	Peter Young, Alice Lai, Micah Hodosh, and Julia Hock-	732
680	Gray, Chelsea Voss, Alec Radford, Mark Chen, and	enmaier. 2014. From image descriptions to visual	733
681	Ilya Sutskever. 2021. Zero-shot text-to-image gener-	denotations: New similarity metrics for semantic in-	734
682	ation. In <i>International Conference on Machine</i>	ference over event descriptions. <i>Transactions of the</i>	735
683	<i>Learning</i> , pages 8821–8831. PMLR.	<i>association for computational linguistics</i> , 2:67–78.	736
684	Robin Rombach, Andreas Blattmann, Dominik Lorenz,	Xuandong Zhao, Kexun Zhang, Zihao Su, Saastha	737
685	Patrick Esser, and Björn Ommer. 2022. High-	Vasan, Ilya Grishchenko, Christopher Kruegel, Gio-	738
686	resolution image synthesis with latent diffusion mod-	vanni Vigna, Yu-Xiang Wang, and Lei Li. 2025. In-	739
687	els. In <i>Proceedings of the IEEE/CVF conference</i>	visible image watermarks are provably removable	740
688	<i>on computer vision and pattern recognition</i> , pages	using generative ai. <i>Advances in Neural Information</i>	741
689	10684–10695.	<i>Processing Systems</i> , 37:8643–8672.	742
690	Chitwan Saharia, William Chan, Saurabh Saxena,	Mingjian Zhu, Hanting Chen, Mouxiao Huang, Wei Li,	743
691	Lala Li, Jay Whang, Emily L Denton, Kam-	Hailin Hu, Jie Hu, and Yunhe Wang. 2023. Gendet:	744
692	yar Ghasemipour, Raphael Gontijo Lopes, Burcu	Towards good generalizations for ai-generated image	745
693	Karagol Ayan, Tim Salimans, and 1 others. 2022.	detection. <i>arXiv preprint arXiv:2312.08880</i> .	746
694	Photorealistic text-to-image diffusion models with		
695	deep language understanding. <i>Advances in neural</i>		
696	<i>information processing systems</i> , 35:36479–36494.		
697	Zeyang Sha, Zheng Li, Ning Yu, and Yang Zhang. 2023.		
698	De-fake: Detection and attribution of fake images		
699	generated by text-to-image generation models. In		
700	<i>Proceedings of the 2023 ACM SIGSAC Conference</i>		
701	<i>on Computer and Communications Security</i> , pages		
702	3418–3432.		
703	Chuangchuang Tan, Renshuai Tao, Huan Liu, Guanghua		
704	Gu, Baoyuan Wu, Yao Zhao, and Yunchao Wei. 2024.		
705	C2p-clip: Injecting category common prompt in		
706	clip to enhance generalization in deepfake detection.		
707	<i>arXiv preprint arXiv:2408.09647</i> .		
708	Chuangchuang Tan, Yao Zhao, Shikui Wei, Guanghua		
709	Gu, Ping Liu, and Yunchao Wei. 2023. Rethinking		
710	the up-sampling operations in cnn-based generative		
711	network for generalizable deepfake detection. <i>arXiv</i>		
712	<i>preprint arXiv:2312.10461</i> .		
713	Paula Dootson T.J. Thomson, Daniel Angus. 2020. 3.2		
714	billion images and 720,000 hours of video are shared		
715	online daily. can you sort real from fake?		
716	Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew		
717	Owens, and Alexei A Efros. 2020. Cnn-generated		
718	images are surprisingly easy to spot... for now. In <i>Pro-</i>		
719	<i>ceedings of the IEEE/CVF conference on computer</i>		
720	<i>vision and pattern recognition</i> , pages 8695–8704.		
721	Zhendong Wang, Jianmin Bao, Wengang Zhou, Weilun		
722	Wang, Hezhen Hu, Hong Chen, and Houqiang Li.		
723	2023. Dire for diffusion-generated image detection.		
724	<i>arXiv preprint arXiv:2303.09295</i> .		
725	H. Wu, J. Zhou, and S. Zhang. 2023. Generalizable		
726	synthetic image detection via language-guided con-		
727	trastive learning. <i>arXiv preprint:2305.13800</i> .		

Appendix

A Dataset Details

A.1 Comparison with Existing AGID Benchmarks

We compare VCT² with two prominent AGID benchmarks:

GenDet (Zhu et al., 2023) introduces a benchmark of roughly 140,000 synthetic images, mostly from COCO and Flickr. It covers several open-source models but lacks prompt diversity from real-world data, and omits proprietary model outputs.

De-Fake (Sha et al., 2023) focuses on detection and attribution using four generative models (DALL-E 2, GLIDE, Latent Diffusion, Stable Diffusion). While well-curated, it lacks coverage of more recent models (e.g., DALL-E 3, Midjourney 6), and does not utilize real prompts from user-facing media.

VCT² provides:

- Over 182,000 images, including 26,000 real image–prompt pairs
- Prompts from both MS COCO and @nytimes (2011–2023)
- Synthetic content from six leading models, including DALL-E 3 and Midjourney 6
- Realistic prompt semantics, verified source integrity, and public accessibility

This makes VCT² a uniquely practical benchmark for evaluating AGID model generalization and real-world robustness.

Table 6 presents 10 records of real images alongside synthetic images generated by different models, based on 10 Twitter prompts for which we had the corresponding Midjourney 6 image (With this model, we generated ~10K synthetic images on MS COCO prompts and only 500 synthetic images on Twitter prompts, as this model blocks image generation for most Twitter prompts.). Full dataset are publically available at [Twitter_{AI}](#) and [COCO_{AI}](#).

B Detection Techniques

This appendix provides detailed descriptions of the 17 AI-generated image detection (AGID) techniques evaluated in our benchmark. These methods span three major detection paradigms: generation artifact-based, feature representation-based, and hybrid approaches. This taxonomy is designed to

reflect the breadth of design assumptions across the literature and serves as the foundation for our performance analysis in Section 4.

Artifact-based methods exploit low-level visual artifacts—such as frequency distortions, edge inconsistencies, or upsampling traces—introduced during the image generation process. Feature-based methods, in contrast, analyze semantic-level inconsistencies by leveraging deep neural representations from CNNs, vision transformers, or CLIP encoders. Hybrid methods combine both low-level and high-level signals, often incorporating alignment objectives or learned fusion strategies to improve robustness.

The detectors described here were selected based on recency, diversity, and public availability, and represent both classical and state-of-the-art AGID strategies. Each technique is evaluated under zero-shot settings using default public checkpoints, and grouped by detection paradigm below.

B.1 Generation Artifact-Based Detection

Generation artifact-based detection techniques focus on identifying visual artifacts produced during the generation process, analyzing both spatial and frequency domains.

Tan (Tan et al., 2023) found that the up-sampling operator introduces artifacts not only in frequency patterns but also in pixel arrangements within images. The authors introduce the concept of Neighboring Pixel Relationships to capture and characterize these generalized structural artifacts caused by up-sampling operations.

Corvi (Corvi et al., 2023) observed that synthetic images, especially those generated by diffusion models like GLIDE and Stable Diffusion, exhibit distinctive differences in mid-to-high frequency signals compared to real images. However, this distinction is less pronounced in images produced by newer models, such as DALL-E and ADM. Although their method accurately distinguishes synthetic and real images in controlled settings, it struggles in real-world scenarios.

Doloriel (Doloriel and Cheung, 2024) explored masked image modeling for universal fake image detection. Their approach involves both spatial and frequency domain masking, leading to a deepfake detector based on frequency masking.

Chen (Chen et al.) enhance detector generalization diffusion generated images by generating hard samples through high-quality diffusion reconstruction. These reconstructed images, which closely

resemble real ones but retain subtle artifacts, train detectors to differentiate between real and generated images, including those from unseen diffusion models.

B.2 Feature Representation-Based Detection

Feature representation-based detection methods distinguish real images from synthesized images by leveraging deep learning models to extract and analyze complex visual features.

Wang (Wang et al., 2020) proposed a universal detector using a ResNet-50 classifier (He et al., 2016) with random blur and JPEG compression data augmentation. When trained on images generated by a single CNN generator (ProGAN), their model demonstrated strong generalization across unseen architectures, including StyleGAN2 (Karras et al., 2020) and StyleGAN3 (Karras et al., 2021).

Mandelli (Mandelli et al., 2022) developed a GAN-generated image detector based on an ensemble of CNNs. Their method emphasizes generalization by ensuring orthogonal results from CNNs and prioritizing original images during testing.

Wang (Wang et al., 2023) introduced a technique that measures the error between an input image and its reconstructed counterpart generated by a pre-trained diffusion model. They observed that diffusion-generated images are more accurately reconstructed than real images, highlighting a key discrepancy for detection.

Wu (Wu et al., 2023) employed language-guided contrastive learning to capture inherent differences in the distributions of real and synthetic images. Their method augments training images with designed textual labels, enabling joint image-text contrastive learning for forensic feature extraction.

Sha (Sha et al., 2023) addressed the challenges of fake image detection and attribution. Their approach involves: (i) building a machine learning classifier to detect fake images generated by various text-to-image models, including DALL-E 2, Stable Diffusion, GLIDE, and Latent Diffusion, and benchmark prompt-image datasets such as MS COCO and Flickr30k (Young et al., 2014); (ii) attributing fake images to their respective generative models to enhance accountability; and (iii) examining how prompts influence detection and attribution performance.

Aghasanli (Aghasanli et al., 2023) introduced a deepfake detection method that combines fine-

tuned Vision Transformers (ViTs) with Support Vector Machines (SVMs). Their method provides interpretability by analyzing the SVMs’ support vectors to distinguish between real and fake images generated by various diffusion models.

Chen (Chen et al., 2024) proposed a straightforward method that extracts the simplest patch from an image and sends its noise pattern to a binary classifier, demonstrating effectiveness with minimal complexity.

Koutlis (Koutlis and Papadopoulos, 2024) utilized intermediate outputs from CLIP’s image encoder for enhanced AI-generated image detection. They introduced a Trainable Importance Estimator to dynamically assess the contributions of each Transformer block, boosting generalization across generative models.

Liu (Liu et al.) presented OCC-CLIP, a CLIP-based framework for few-shot one-class classification. This method is particularly effective when only a few images generated by a model are available, and access to the model’s parameters is restricted. OCC-CLIP combines high-level and adversarial data augmentation techniques to attribute images to specific generative models accurately.

To enhance generalization to unseen generative models, Ojha (Ojha et al., 2023) propose a novel approach that avoids explicitly training a classifier to distinguish real from fake images. Instead, their method leverages the feature space of large pre-trained vision-language models and employs techniques such as nearest neighbor classification.

Tan (Tan et al., 2024) enhance the image encoder’s ability to detect deepfakes by integrating category-related prompts into the text encoder of CLIP.

B.3 Hybrid Techniques

Hybrid techniques combine low-level artifact analysis with high-level semantic feature extraction to effectively distinguish AI-generated images from real ones.

Yan (Yan et al., 2024) propose a hybrid-feature model that integrates high-level semantic information (using CLIP) with low-level artifact analysis to improve detection robustness.

Liu (Liu et al., 2024) incorporate a forgery-aware adapter that integrates local forgery traces from both image and frequency domains. Their method employs language-guided alignment, using contrastive objectives between image features and text prompts to enhance generalization.

To guide our benchmark evaluation, we selected 17 state-of-the-art AGID methods spanning all three categories. This categorization enables us to evaluate model robustness from complementary perspectives: from low-level artifact exploitation to high-level semantic inconsistency analysis. This taxonomy is used to analyze generalization performance in Section 5.

C Detection Performance Overview

Tables 2 and 3 provide an overview of the performance of different detection techniques across synthetic datasets generated from MS COCO and Twitter prompts, respectively. The metrics measured are Accuracy (Acc), Recall (R), and Precision (P), providing insights into each model’s ability to differentiate real from AI-generated images.

C.1 Performance by Detection Technique

- **CNNDetection, NPR and Fake Image Detection:** These methods showed variable results, characterized by low recall but higher precision across several models. This indicates a tendency to correctly identify generated images when detected, but with many instances being missed (false negatives).
- **DM Image Detection and De-Fake:** DM Image Detection demonstrated high precision across all models, particularly excelling with Stable Diffusion versions and Midjourney 6, effectively capturing generated images. De-Fake consistently maintains strong metrics across SD (2.1, XL and 3), DALL.E 3, and Midjourney 6 but struggles with SD 3.5 Large images, exhibiting lower accuracy, precision, and recall. This drop in performance likely results from SD 3.5’s refined generation and post-processing that minimize the artifacts and noise patterns AGID techniques depend on.
- **GAN Image Detection, SSP and DIRE:** These methods had mixed performance, particularly excelling in precision.
- **DRCT (ConvB and UnivB):** Both versions of DRCT showed strong accuracy, recall, and precision across most models but experienced a slight performance drop with Midjourney 6, indicating challenges with proprietary models.

- **OCC-CLIP and Deep Fake Detection:** OCC-CLIP had lower recall with SDXL but balanced performance for DALL.E 3 and Midjourney 6; while Deep Fake Detection demonstrated steady, consistent performance, with all of its metrics remaining within a similar range.
- **Universal Fake Detect:** Universal Fake Detect performed better on SD (2.1, XL, and 3) models but its performance dropped when applied to SD 3.5, DALL.E 3, and Midjourney 6. Notably, we observed a significant increase in recall for DALL.E 3-generated images across both datasets.
- **C2P CLIP:** C2P CLIP consistently performs poorly with low accuracy and recall, clearly showing that it often misses AI-generated images. Although its precision remains high across both datasets overall, it declines significantly for Midjourney 6 images in both datasets and for SD3 images in the Twitter dataset.

The results indicate that there is no one-size-fits-all solution for detecting AI-generated images. Different generative models pose unique challenges, and the performance of each detection method varies based on its ability to identify specific artifacts. De-Fake and DRCT (ConvB and UnivB) were the most consistent performers, highlighting their robustness across models. Future research should aim to improve detection for proprietary models like SD 3.5 Large, Midjourney 6 and DALL.E 3, where many techniques struggled.

D Visual AI Index Overview

Midjourney 6 achieved the highest Visual AI Index (V_{AI}) score on MS COCO prompts, indicating superior visual coherence and quality compared to other models. Stable Diffusion 2.1 also showed relatively high performance, suggesting that diffusion-based methods can achieve strong visual results but may still be outperformed by proprietary methods like Midjourney 6.

From the V_{AI} scores on Twitter prompts, Midjourney 6 remained the top performer, followed closely by DALL.E 3. This suggests that proprietary models are particularly robust in generating high-quality images, even with more diverse and potentially less structured prompts. The accuracy heat maps in Figure 5 also highlight differences in

how AGID methods perform across models. Methods like De-Fake and DRCT were particularly effective at detecting Midjourney-generated images, whereas detection on DALL.E 3 and SDXL proved more challenging. This indicates that the texture and artifact characteristics differ significantly across these models, affecting detection reliability. These results underscore the challenges faced by AGID methods when applied to high-quality proprietary models. While some detection methods, like De-Fake and DRCT, performed consistently well, the V_{AI} scores reveal that generated image quality plays a significant role in detection difficulty. Future work should focus on improving the robustness of detection techniques against models that prioritize high visual fidelity, such as Midjourney 6 and DALL.E 3.

E Supplementary Figures

The supplementary figures included in this appendix (Figure 13 and Figure 12) provide additional insights and visual examples that support and enhance the concepts and results discussed in the main paper.

F Visual AI Index (V_{AI})

F.1 Texture and Semantic Cues

Feature-space PCA plots and LBP texture visualizations (Figures 6) highlight differences in object coherence, sharpness, and semantic consistency. Newer generators often produce globally coherent scenes that are statistically inconsistent with real-world image distributions, particularly in terms of texture smoothness and semantic density. These deviations offer forensic cues that detectors can exploit.

F.2 LBP Texture Analysis

Local Binary Pattern (LBP) is commonly used for texture analysis, image recognition, and quality assessment. LBP plots can indirectly assess image quality, as sharper images generally produce more distinct patterns in their LBP representations. If the LBP pattern appears blurred or lacks clear edges, it may indicate a loss of detail or lower resolution in the image. AI-generated images sometimes lose fine-grained texture, which would be visible as less distinctive LBP features. In Figure 6 we can see that image generated by Midjourney has specific facial textures and subtle expression lines whereas image generated by SD 3 has inconsistencies and

lack of texture in certain areas. facial features, facial structures, hair lines, edges in clothing, and wrinkles are preserved in each segment for the Midjourney image but SD 3 image completely lost it.

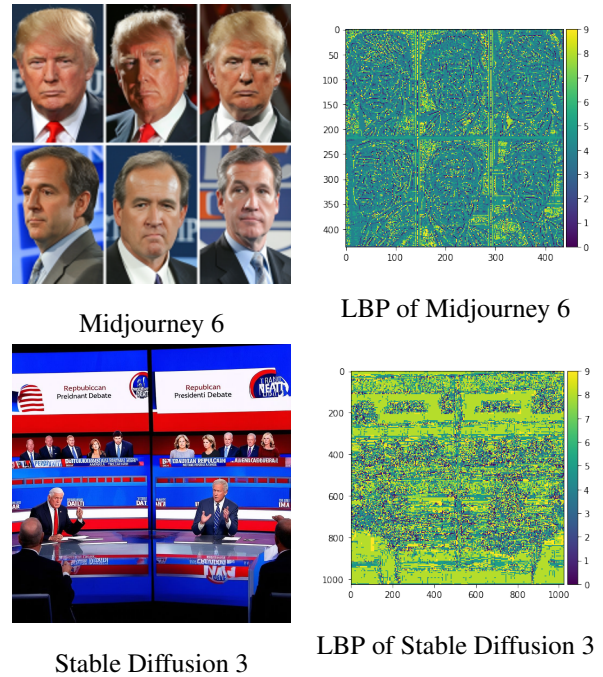


Figure 6: Comparative analysis of texture patterns in images generated by different T2I models using Local Binary Pattern (LBP) representation.

F.3 Pairwise Scatter Plot Analysis

The pairwise scatter plots in Figure 7, 8, 9, 10, and 11 reveal clear differences in model distributions. DALL.E 3 and SDXL maintain higher object coherence across varying texture complexities, indicating stronger object integrity. In contrast, SD 2.1 and SD 3 show more dispersed patterns, reflecting lower consistency. These trends underscore evolutionary gains, with newer models like SD 3 and SDXL offering improved color distribution and contextual relevance. Midjourney 6 and DALL.E 3 form well-separated clusters, suggesting better handling of coherence and complexity. Stable Diffusion variants display mixed behaviors, revealing inconsistencies in texture and boundaries; consistent with their lower V_{AI} scores. Overall, the analysis highlights the importance of balancing texture complexity and object coherence to improve generative and detection model performance.

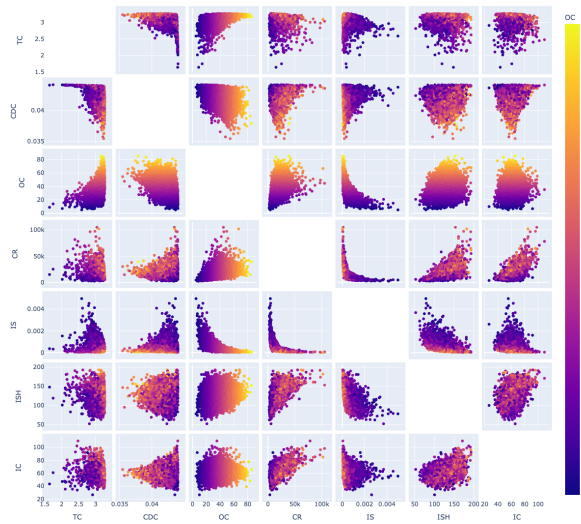


Figure 7: DALL-E 3

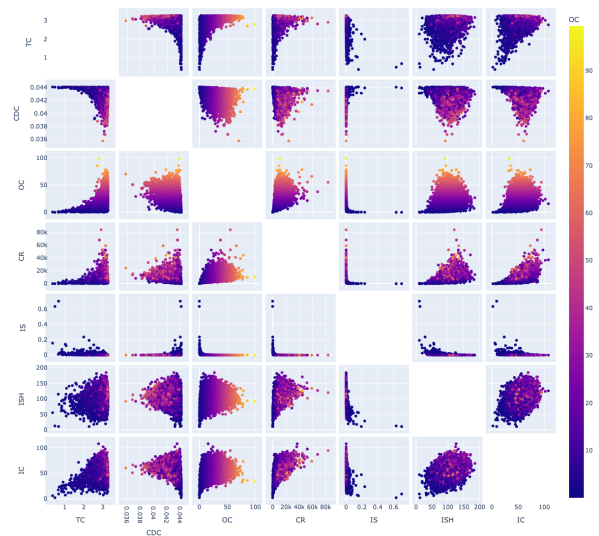


Figure 10: Stable Diffusion 2.1

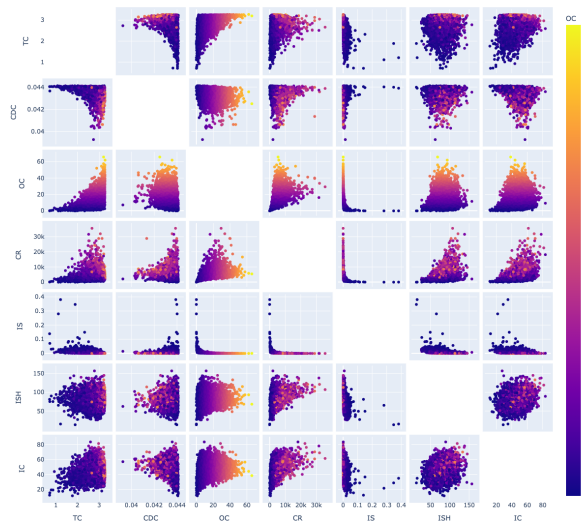


Figure 8: Midjourney 6

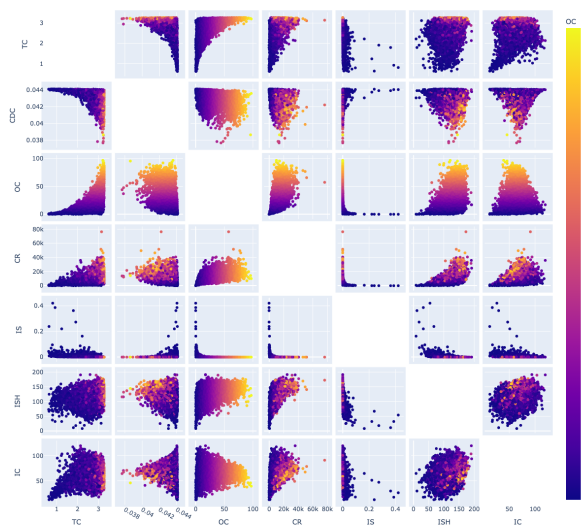


Figure 9: Stable Diffusion 3

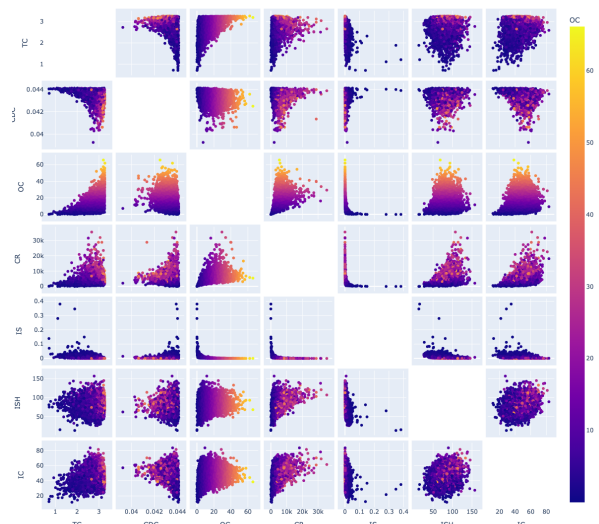


Figure 11: Stable Diffusion XL



DALL·E 3



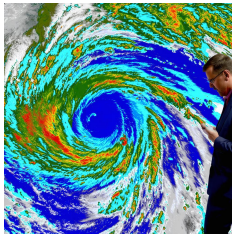
Midjourney 6



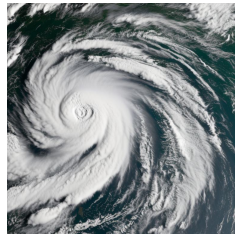
DALL·E 3



Midjourney 6



Stable Diffusion 3



Stable Diffusion 2.1



Stable Diffusion 3



Stable Diffusion 2.1



Stable Diffusion XL



Stable Diffusion XL

Figure 12: Generated images by different AI models for the prompt: "Have you ever wondered why we name hurricanes? The New York Times meteorologist Judson Jones explains."

Figure 13: Generated images by different AI models for the prompt: "At least six candidates appear to have made the cut so far for the second Republican presidential debate on Sept 27. See which candidates have and have not qualified so far."

Table 6: Real images and synthetic images generated by different models.

Real Image	SD2.1	SDXL	SD3 Medium	SD3.5 Large	DALL.E 3	Midjourney 6
