# Towards Explainable Fake Image Detection with Multi-Modal Large Language Models

Yikun Ji
da-kun@sjtu.edu.cn
Shanghai Jiao Tong University
Shanghai, China

Yan Hong
Ant Group
Hangzhou, China

Jiahui Zhan
Shanghai Jiao Tong University
Shanghai, China

Haoxing Chen
Ant Group
Hangzhou, China

Jun Lan
Ant Group
Hangzhou, China

Huijia Zhu*
Ant Group
Hangzhou, China

Weiqiang Wang
Ant Group
Hangzhou, China

Liqing Zhang*
Shanghai Jiao Tong University
Shanghai, China

Jianfu Zhang*
c.sis@sjtu.edu.cn
Shanghai Jiao Tong University
Shanghai, China

## Abstract

Recent advancements in image generation have provoked social and security concerns, yet most detection methods rely on black-box models that generalize poorly. By utilizing advances in Multi-modal Large Language Models (MLLMs), we propose a framework that fuses six specialized paradigms, each analyzing a distinct aspect of the image, to provide a final verdict with coherent, evidence-based reasoning. Experimental results on a diverse dataset of real and AI-generated images demonstrate that our approach outperforms both traditional detection methods and top humans, while providing explainability. This study underscores the potential of MLLMs in developing robust, explainable, and reasoning-driven detection systems. The code is available at https://github.com/Gennadiyev/mllm-defake.

## CCS Concepts

• **Computing methodologies** → **Natural language generation**; *Computer vision*; • **Social and professional topics** → *Censorship*.

## Keywords

Generative AI, Multi-modal Large Language Models, Explainable AI, AI safety

## 1 Introduction

The rapid emergence of deep learning and artificial intelligence, exemplified by Generative Adversarial Networks (GANs) [22] and diffusion models [12], has revolutionized synthetic image generation, creating remarkable opportunities for creative innovation. Yet, these technologies also pose critical societal challenges, particularly in cyber-security, where the proliferation of synthetic images can threaten both individual and public safety [35].

Detecting AI-generated images remains a challenging task, usually formulated as a binary classification problem. While deep learning-based methods have exhibited solid performance [64, 71], they are prone to overfitting and often rely heavily on large, labeled datasets of real and synthetic images. Alternative approaches

*Corresponding Authors



**Figure 1: Traditional methods do not provide reasons for their predictions, while our method provides sensible reasons behind the verdict.**

focus on statistical anomalies [3, 56] but frequently struggle with poor generalization, especially when limited to specific generative models or image domains (*e.g.*, facial images). Consequently, there is a pressing need for detection models that not only perform well across diverse generative techniques but also provide interpretability, enabling users to understand why images are flagged.

With the development of multi-modal large language models (MLLMs), researchers have begun exploring their potential in detecting AI-generated images. A key advantage of using MLLMs is their ability to generate predictions with explanations, thus enhancing interpretability. Studies have shown that an open-source MLLM (mPlug-owl2 [75]) with 7B parameters can achieve an accuracy of 71.78%, while OpenAI's GPT-4V reaches 78.03%, outperforming the average human accuracy of 74.51% [38]. Zhang et al. [83] highlights the quality of reasons that MLLMs provide when verifying image authenticity. Although these models demonstrate notable potential, MLLMs remain susceptible to hallucination [19] and generally cannot surpass the accuracy of traditional methods. Meanwhile, their responses can be sensitive to prompt design [72], suggesting the need for more refined techniques to enhance both accuracy and reasoning quality.

In this work, we enhance the capabilities of multi-modal large language models (MLLMs) through an innovative interrogation methodology. By integrating human expertise with MLLM-based reasoning, we design six distinct prompts, each focusing on a distinct aspect of the image. By fusing the model responses from these

**P1 (Defect Query):**

···*features of AI-generated / real images?*

{Features of AI-generated / real images}

Sample A  Sample B  *real or generated···?*

[Sample A]…shadows are correct…no evident stitching errors…background is properly blurred…**Likely real.**
[Sample B]…reflections on the diamonds appear slightly exaggerated…unrealistic reflections and perfect surfaces…**Generated.**

**P2 (Regional Analysis):**

Sample A  Sample B  ···*focus on the main object as* ROI A  ROI B *provided*···

[Sample A]…highly reflective pupils…fur lacks depth…unnatural textures… **AI-generated.**
[Sample B]…diamond reflections behave realistically…skin texture…highly detailed, with visible pores and nail ridges… **Likely real.**

**P3 (Common Sense Reasoning):**

···*analyze* Sample A  Sample B ···*Consider:*
1. *Physical proportions*···
2. *Spatial relationships*···
3. *Natural world physics and rules*···
4. *Other logical inconsistencies*···

[Sample A]…dog's anatomy appears correct…no extra limbs…lighting exhibits hyper-realistic characteristics…lacks subtle imperfections…**Likely generated.**
[Sample B]…diamond reflections and metallic finish appear plausible… text appears slightly out of place… lighting and shadows align correctly…**Likely real.**

⚙ System  User  ✗ Assistant  ✂ Verdict

**P4 (Few-Shot):**

Exemplar A  Exemplar B ···*real or generated? Why?*
{preset reason}

Sample A  Sample B  *real or generated···?*

[Sample A]…a cartoon-like dog…overly smooth fur, overly defined reflections, and uniform patterns…**AI-generated.**
[Sample B]…maintains an appropriate balance of realism, imperfections, and natural variation…retains natural imperfections…**likely real.**

**P5 (Structural Analysis):**

Sample A  Sample B  ···*class_name?*
{class_name}  **Payload 1**

···*list the key components of {class_name}*···?
{components}  **Payload 2**

Sample A  Sample B  ···*shows a {class_name}· Show me a list of its key components*···
{components}

···*examine for these structural inconsistencies:*
1. *Repetitive patterns that are too uniform*
2. *Missing essential components*
3. *Incorrect object placement*  **Payload 3**

[Sample A]…No missing facial features…No misplaced shadows…background is clearly not on the same layer…**could be real.**
[Sample B]…ring edges are unnaturally smooth…diamonds are placed with mechanical precision… too perfect…background shows repetitive pattern…**Likely generated.**

**P6 (Stereotype Matching):**

Sample A  Sample B  ···*main subject depicted?*
{class_name}  **Payload 1**

···*what are common stereotypes or patterns of {class_name}*···?
{stereotype_response}  **Payload 2**

Sample A  Sample B  ···*shows {class_name}··· has the following stereotypes: {stereotype_response}···Examine for these attributes*···  **Payload 3**

[Sample A]…symmetry in the dog's eyes and face is strikingly perfect…fur exhibits uniformity…depth of field and background blurring seem more **like AI render**…
[Sample B]…diamonds exhibit strong sparkle…aligns with expected optical behavior…metal band has smooth but not overly uniform texture…**Likely real.**

**Fusion:**

···*Use the previous responses effectively, combine the information, and make a decision*···

···*make a decision based on the information provided*···

[Sample A]All previous analysis indicate that the image is likely generated rather than real. The strongest evidence comes from... **Final Conclusion: AI-Generated**
[Sample B]After evaluating the image using all six paradigms, there are no clear signs of AI generation… **Final Decision: Real Image**

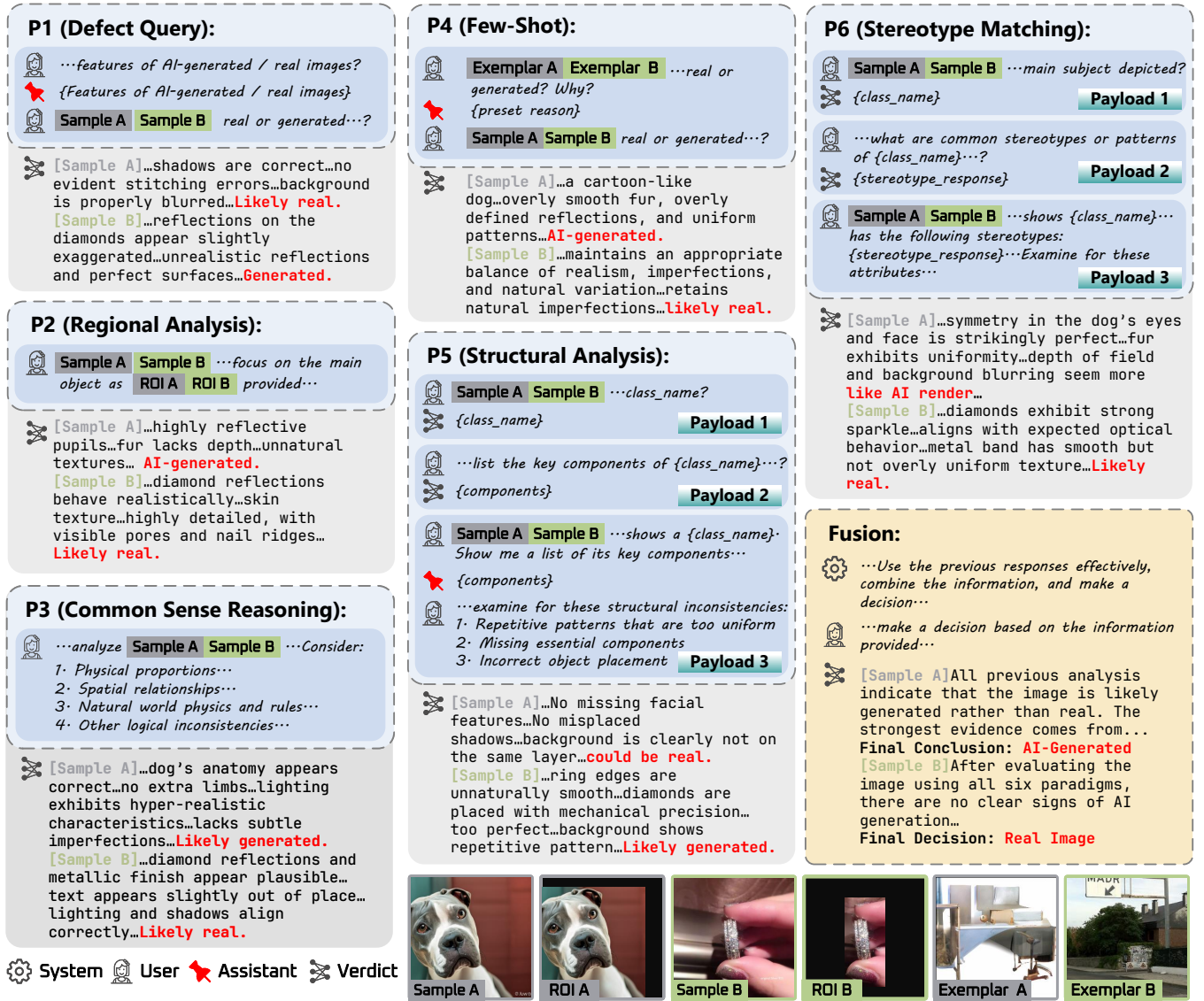Sample A  ROI A  Sample B  ROI B  Exemplar A  Exemplar B

**Figure 2: The overall design of the proposed MLLM-based AI-generated image detection framework.**

prompts, we achieve a higher overall accuracy and improved interpretability. Our quantitative and qualitative analysis reveals that MLLMs can process visual information similarly to human perception when distinguishing between real and AI-generated images. The models exhibit strong generalization capabilities, with GPT-4o achieving 93.4% accuracy, outperforming traditional deep learning models [71] (91.8%) and AEROBLADE [56] (85.2%), as well as the most accurate human annotator (86.3%).

This framework highlights the promise of MLLMs for robust, transparent, and explainable AI-generated image detection. Moreover, the model can explicitly articulate why an image is classified as generated, a crucial capability for platforms aiming to provide clear explanations to users, as illustrated in Figure 1.

## 2 Related Works

### 2.1 AI-Generated Image Detection Methods

AI-generated image detection has received significant attention due to the rapidly improving fidelity of synthetic images produced by Generative Adversarial Networks (GANs) [22], autoregressive transformers[68], and diffusion-based models [15, 27, 57, 62]. Early research primarily focused on artifact detection, extracting discriminative cues from spatial and frequency domains to highlight inconsistencies between real and synthetic images [5, 14, 17, 44, 81]. These inconsistencies often stem from upsampling artifacts in GAN pipelines [17], texture synthesis mismatches[44], or limited high-frequency decay in generated outputs[14]. By learning to detect these subtle traces, models can effectively expose fakes, even when artifacts are visually imperceptible.

Beyond artifact-based methods, deep-learning based methods are also advancing. ResNet and Vision Transformers trained on real and synthetic samples can effectively distinguish real from AI-generated images [53, 64, 71]. This paradigm leverages deep networks' strong feature extraction capabilities to automatically learn discriminative patterns. However, the effectiveness of large-scale training heavily depends on data diversity and quality. Generalization to unseen generative models, particularly from GANs to diffusion-based methods, remains challenging. Meanwhile, as diffusion models produce increasingly photorealistic images, traditional artifact-based cues, including unnatural hues, faulty perspectives, and blurry text, are no longer reliable. Additionally, Grommelt et al. [23] highlighted that existing models are prone to overfitting on JPEG artifacts, further limiting their robustness.

One branch of complementary research direction focuses on improving model explainability. Most existing detection methods output only a binary classification, failing to indicate how or where synthetic cues are identified. Recent efforts emphasize fine-grained or localized detection, providing more transparent classification rationales. Notable approaches include measuring reconstruction errors between an input image and its diffusion-based reconstruction [56], leveraging multi-branch systems to produce multi-level labels for generated images [3], or computing local intrinsic dimensionalities [45]. Some approaches integrate text-image contrastive learning for better explainability [60, 61, 63].

While these methods provide regional clues, they lack human-understandable explanations, limiting their effectiveness in enhancing explainability.

## 2.2 Multi-Modal Large Language Models

Multi-modal large language models (MLLMs) are large language models (LLMs) enhanced with vision or audio modalities [2, 40]. Over the past few years, both open-source and proprietary models have significantly advanced vision-language capabilities. Open-source ones include Qwen-VL [1, 55, 70] and QVQ [54], InternVL [7], Ovis [47], LLaVA and its predecessors [41–43], Llama-3 [67], among many others. Proprietary models have also progressed rapidly. Notable releases include GPT-4o, GPT-4o-mini [50, 51] and Gemini-1.5 [20, 21]. These models demonstrate remarkable performance in vision-related tasks, including image reasoning and interpretation. Most benchmarks for vision-language models focus on reasoning capabilities related to the image [39, 77] or geometric problem-solving skills [31, 46].

Li et al. [38] proposed FakeBench, a benchmark for general AI-generated image detection that includes three primary datasets: reasoning, interpretation, and open-question answering. It covers ten generative models across diverse image content, revealing that current MLLMs demonstrate moderate detection accuracy, basic reasoning abilities, and limited fine-grained authenticity analysis. Zhang et al. [83] created A-Bench, which "challenges both high-level semantic understanding and low-level visual quality perception" of MLLMs in the form of multiple-choice questions, pointing out major flaws in each AI-generated image. The question-answer pairs are curated by humans and not actively generated by MLLMs. FakeReasoning [18] conducts the forgery reasoning task in structured steps using MLLMs, giving promising results by fine-tuning

LLaVA-1.5-13b on the MMFR-Dataset, created using GPT-4o and human experts.

Our proposed method significantly enhances the accuracy while also providing fine-grained explanations of why an image is considered real or generated, addressing these limitations.

## 3 Methodologies

Given an image $x$, we employ a detection function $f$ to classify it as either real (*i.e.*, natural) or fake (*i.e.*, AI-generated), producing a prediction $y = f(x)$. Our approach is driven by the observation that different prompt formulations enable MLLMs to focus on distinct aspects of an image.

Techniques such as image cropping [79], few-shot learning [25, 84], and region-of-interest selection [73] enhance MLLM performance by refining visual feature extraction and interpretation. These techniques influence how models extract and analyze relevant image attributes, leading to more effective detection. Prompt formulation plays a critical role in MLLM performance. We design optimized prompts that effectively integrate textual and visual information, mimicking human perception. By combining domain expertise with MLLM reasoning, our approach enhances both explainability and accuracy. The proposed framework is illustrated in Figure 2, and the following sections provide a detailed explanation of our prompt construction and detection pipeline.

### 3.1 Prompt Engineering

We provide visual images and textual prompts to the MLLMs, leveraging chain-of-thought (CoT) prompting to improve reasoning. Our approach structures MLLM *sessions* (*i.e.*, context window) as sequential queries, where each query is defined as: $o_i = M(t_i, v_i, c_{i-1})$, where $i$ denotes the current query, $M$ represents the MLLM, $o$ is the model's output, $t$ and $v$ are text/visual prompts (image type in this task) from the user, and $c$ is the session context (query history) from previous queries. Each query triggers a computational *payload*, incurring processing costs. The output $o$ is incorporated into the next query as part of the updated context: $c_{i+1} = (t_i, v_i, o_i, c_i)$. To ensure consistency and mitigate randomness, we utilize a predefined $o$ as an *assistant* reference, summarized by the user. This approach also optimizes inference efficiency by reusing previous outputs, improving model stability and reducing computational costs while maintaining fast inference speed.

*Verdict Prompt (P0).* In our method, a designated *verdict prompt* $n$ is introduced, where the input image $v_n$ corresponds to $x$, and the prompt $t_n$ is structured as follows:

```
[Is this image real or fake? End your response with
either "real" or "fake".]
```

This enforces a concise model prediction, allowing us to directly use the output $o_n$ as the final classification result $y$. P0 serves as the baseline of our framework, leveraging MLLMs' classification capabilities. Unless otherwise specified, our approach always concludes with P0.

To systematically explore prompt variations, we develop six distinct prompt strategies, categorized into three types: General Prompts (direct analysis mimicking human reasoning), Few-Shot Prompts (providing a few examples to guide model predictions) and Content-Based Prompts (first analyzing image content before

predicting). We briefly introduce each prompt type later. For full implementation details, please refer to the provided code.

*System Prompt.* To enhance the interoperability of our method, we define an initial system prompt $h_0$ that defines the session's general settings:

```
[You are an AIGC detection specialist. The user will
ask you about whether an image is real or generated.
Observe the image carefully to decide whether it is
real or generated. Explain your reasons, and end your
response with either "real" or "generated".]
```

This ensures that the model justifies its classification decision, maintaining explainability in the prediction function $f$. For all subsequent prompts P1–P6, this system prompt initializes the session.

## 3.2 General Prompts

Humans analyze images top-down, first grasping the overall scene before focusing on key objects and details. We design our model to follow a similar process, leveraging MLLMs' reasoning abilities to detect anomalies. Our framework systematically analyzes primary objects, assessing their properties against real-world constraints and logical coherence. We design three prompts: Defect Query (P1): Conducts a broad inspection to identify general anomalies; (P2): Focuses on specific regions to refine local anomaly detection; and Common Sense Reasoning (P3): Applies logical reasoning to verify plausibility against real-world knowledge. By integrating these perspectives, our approach enhances the reliability and explainability of AI-generated image detection.

*Defect Query (P1).* Generative models rely on extensive real-world data, yet no dataset can fully encapsulate all causal factors of a scene. Photography follows a cause-and-effect paradigm—objects are shaped by lighting, camera angles, lens parameters, and the photographer's intent. Models trained solely on photographic data often fail to internalize these relationships, leading to artifacts such as incorrect shadows, unnatural object placements, and abnormal focus or blur. To address this, the defect query prompt directs the MLLM to identify such anomalies. The model is informed that the image may contain defects and is tasked with listing potential issues. The process consists of three queries. In the first two queries, $o_1$ lists known AI-generated artifacts, $o_2$ highlights features typical of real images, and $t_1$ / $t_2$ are

```
[What are common defects in AI-generated images?] /
[What are the features that real images often have?]
```

Finally the model makes a final prediction with P0, reasoning based on $o_1$ and $o_2$. This query structure enables MLLMs to predict and justify their classification by identifying defects such as abnormal object proportions, unrealistic reflections, etc.

*Regional Analysis (P2).* Each input image contributes a fixed number of tokens to the MLLM. To emphasize key areas, this method extracts Regions of Interest (ROI) from the original image, allowing the model to focus on specific objects or sub-scenes, such as facial features or text on a book. While this approach may omit broader context, it enhances detail analysis in a controlled setting—especially useful when the main object occupies only a small portion of the frame. This method is inspired by prior work on synthetic image detection Simonyan et al. [61], where heatmaps are

used to improve model explainability. The process consists of two queries. In the first query, $v_1 = x$, $o_1$ will request user confirmation of the ROI, and $t_1$ is

```
[Instead of analyzing the full image, focus on the
ROI in the next image to determine whether it is real
or generated.]
```

We extract the top three most relevant regions using DINOv2 [52], generating heatmaps that are thresholded to form bounding boxes from binary pixel maps. Pretrained weights from ImageNet are used. For the second query, $v_2$ is the ROI image and $t_2$ is

```
[Here is the ROI. Please indicate whether the image
is real or generated based on ROI. End your response
with either "real" or "generated"].
```

The final prediction is then generated, along with an explanation focusing on the main area of the image. This ensures that the model prioritizes fine details of key objects, such as pattern distortions, unrealistic textures, and other detail problems, enhancing detection accuracy.

*Common Sense Reasoning (P3).* AI-generated images often contain logical inconsistencies that are implausible in the real world, such as incorrect numbers of fingers, unnatural object attachments, and distorted text. While humans easily recognize these anomalies, MLLMs may overlook them when prioritizing overall image realism. To address this, and inspired by [82], we prompt the MLLM to detect inconsistencies by comparing images against real-world patterns. This method uses a single query with $v_1 = x$ and $t_1$ is

```
[ Please analyze this image for any violations of
common sense or real-world logic. Consider: 1. Physical
proportions (e.g., number of fingers, limbs); 2. Spatial
relationships between objects, such as merged textures;
3. Natural world physics and rules, such as unnaturally
bent fences; 4. Other logical inconsistencies in the
scene Explain any inconsistencies you find. If you find
many unexplainable inconsistencies, the image is likely
generated.]
```

This approach guides the model to examine critical visual cues, such as limb counts and spatial relationships. While these anomalies may be subtle, they serve as strong indicators of AI-generated content. By incorporating logical consistency checks, this method improves model generalization across diverse images while enhancing interpretability in "generated" classification.

## 3.3 Few-Shot Prompt

*Few-Shot (P4).* In the realm of MLLM, few-shot learning refers to the technique of providing exemplars in the prompt before the query. Our few-shot prompt provides two labeled exemplars along with their human-annotated responses to a baseline prompt, refining the model's conceptual understanding. Our prompt design consists of three turns, each using a structured text prompt similar to P0 for prediction. The first two queries introduce labeled exemplars, one real for $v_1$ and one fake for $v_2$, with manually annotated justifications for their classification. These responses replace $o_1$ and $o_2$, and are packaged with previous images in a single history payload. Finally, P0 will be conducted for the final prediction. This method leverages the MLLM's reasoning and generalization capabilities under data constraints, ensuring interpretability. By

referencing only a few exemplars, the model's decision process remains transparent and grounded in direct comparisons. The compact dataset encourages the model to focus on essential features, such as object outlines, textures, and logical color consistency, instead of relying on memorized distributions, leading to more robust conclusions.

## 3.4 Content-Based Prompts

Content-Based Prompts evaluate an image's internal coherence by analyzing its main subject, structural patterns, and alignment with real-world expectations. These prompts leverage the MLLM's ability to identify key features and provide interpretable justifications for inconsistencies. To prevent conflicts between different queries, the process follows a structured multi-session approach. The first session is Object Identification, where $t$ is

```
[What is the main subject in the image? Please answer
concisely with a short descriptive noun phrase, e.g.,
'a red bus', 'a daisy'.]
```

The second session then builds upon the response from the first turn, refining the analysis with $v = x$. $o$ will be the {class} of $x$. Then the {class} information is integrated into the following sessions for the final prediction, which will be detailed in the next paragraphs.

*Structural Analysis (P5).* This method systematically identifies internal inconsistencies within an image based on its object class. The detection process follows a two-step query approach: 1. Component Identification, where $t$ is

```
[This image shows a {class}. List its key components
before determining whether it is real or generated.]
```

$v = x$, and $o$ outputs a list of expected components based on the identified object class. 2. Structural Verification: The model checks for missing or misplaced components, assessing whether these structural faults indicate AI generation. At the same time, P0 is applied to finalize the classification. This modular approach enhances interpretability by mapping each detected anomaly to a logical explanation, conditioned on the expected class structure. Using this prompt, the model can effectively detect AI-generated images with: unnaturally uniform repetitions (*e.g.*, identical windows in a building, consistent distortions across a fence); absence of critical features (*e.g.*, a dog without a tail, a person missing hands); implausible positioning of objects (*e.g.*, a tree growing out of a building, a car parked on a rooftop).

*Stereotype Matching (P6).* This method first classifies the dominant object in an image, then evaluates whether its features conform too closely or uniformly to common stereotypes. Examples include excessively symmetrical human faces, unnaturally uniform animal textures, or exaggerated stereotypical attributes (*e.g.*, a red-beaked rooster with atypical pigmentation). The detection will first conduct stereotype analyses session, where $t$ is

```
[If an image shows {class}, what are some common
stereotypes or patterns you would expect? List at least
three.]
```

$o$ outputs common stereotypes of the identified class. Then, we start a new session. The model analyzes the image for the presence of these common stereotypes, explaining any anomalies and determining whether the image is real or AI-generated with P0. A certain

degree of stereotypical accuracy enhances realism, but excessive adherence to these patterns may indicate AI generation, reflecting biases learned from training data. By prompting the MLLM to assess how closely an image follows stereotypical features, this method enhances interpretability in detecting unnatural uniformity or exaggerated regularities that deviate from real-world variation.

## 3.5 Fusion Process

Prompts categorized as *General*, *Few-Shot*, and *Content-Based* enhance AI-generated image detection via MLLMs by targeting interpretability aspects like local inspections (P2, P4), logical consistency (P3, P5), and real-world traits (P1, P6). We propose two fusion methods to determine the final verdict. The primary approach executes P1–P6 in parallel before consolidating an output, balancing accuracy and explainability. Alternatively, a majority vote method tallies individual prompt verdicts, trading some explainability for computational efficiency while remaining effective in most cases. Experimental results validate these trade-offs, showing the sequential method excels in interpretability and the majority vote offers practical efficiency.

## 4 Experiments

## 4.1 Experimental Setup

We ensure relevance and effectiveness in real-world AI-generated image detection by covering a diverse range of generation models, including but not limited to GANs, auto-regressive models, and diffusion models. To achieve this, we selected 1,000 real images and 1,000 AI-generated images from WildFake [29], spanning different generation architectures, with detailed distribution provided in the *appendix*. To benchmark MLLM performance against humans, we invited 24 volunteers to classify the images, 12 of whom had prior experience with AI image generation.

We tested recent MLLMs such as GPT-4o, GPT-4o-mini, Llama 3.2 VI (Llama-3.2-Vision-Instruct 11B) [67], LLaVA-CoT [74], Qwen-VL [1, 54, 55, 70], Ovis [47] and InternVL series [7]. Among them, we selected four models for our full analyses: GPT-4o, GPT-4o-mini, Llama 3.2 VI, and LLaVA-CoT. In our experiments, we used OpenAI GPT-4o that points to `gpt-4o- 2024-08-06` and GPT-4o-mini that points to `gpt-4o-mini- 2024-07-18`. The Llama 3.2 VI refers to the Llama-3.2-Vision-Instruct 11B model, from which LLaVA-CoT is a fine-tuned chain-of-thought model. Local MLLMs are deployed using `vllm-0.7.2` [36] on four NVIDIA A100-40G GPUs.

To compare the performance with traditional methods, we select AEROBLADE [56], CNNSpot [71], CommunityForensics [53] and ObjectFormer [69]. We evaluate them on the same sampled dataset used for MLLM and human evaluations. For AEROBLADE, we employ Stable Diffusion versions 1.5, 2.1, and 3.5-large [15, 57] as reconstruction models and select SD3.5-Large for the results, since it provides the best performance. For CNNSpot and Object-Former, since the pretrained models performed suboptimally, we evaluated on checkpoints that we trained on the same data. CommunityForensics requires no training.

Table 2 reports the average inference time per image when six prompts are queried both sequentially and with parallelism enabled. For Llama 3.2 VI [67] and LLaVA-CoT [74], profiling was conducted on 8× NVIDIA A100-40G GPUs. Network latency is included for

| Prompt | GPT-4o | | | GPT-4o-mini | | | Llama 3.2 VI | | | LLaVA-CoT | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | All | Real | Generated | All | Real | Generated | All | Real | Generated | All | Real | Generated |
| P0 | 0.852 | 0.855 | 0.849 | 0.769 | 0.910 | 0.628 | 0.561 | 0.645 | 0.477 | 0.594 | 0.768 | 0.420 |
| P1 | 0.846 | 0.802 | 0.888 | 0.776 | 0.972 | 0.579 | 0.561 | 0.645 | 0.477 | 0.592 | 0.749 | 0.435 |
| P2 | 0.848 | 0.824 | 0.871 | 0.776 | 0.846 | 0.706 | **0.595** | 0.637 | 0.552 | - | - | - |
| P3 | **0.896** | 0.922 | 0.869 | 0.701 | 0.721 | 0.680 | 0.531 | 0.412 | 0.650 | 0.600 | 0.619 | 0.580 |
| P4 | 0.834 | 0.801 | 0.867 | **0.860** | 0.874 | 0.846 | 0.497 | 0.525 | 0.468 | 0.526 | 0.583 | 0.469 |
| P5 | 0.841 | 0.813 | 0.868 | 0.734 | 0.732 | 0.735 | 0.497 | 0.297 | 0.697 | **0.610** | 0.752 | 0.467 |
| P6 | 0.888 | 0.880 | 0.895 | 0.753 | 0.767 | 0.738 | 0.584 | 0.563 | 0.604 | 0.586 | 0.603 | 0.568 |
| w/o P1 | 0.895 | 0.917 | 0.872 | **0.829** | 0.862 | 0.796 | 0.605 | 0.512 | 0.698 | 0.669 | 0.685 | 0.653 |
| w/o P2 | 0.863 | 0.850 | 0.876 | 0.790 | 0.823 | 0.757 | 0.609 | 0.534 | 0.683 | 0.674 | 0.719 | 0.628 |
| w/o P3 | 0.871 | 0.850 | 0.891 | 0.799 | 0.835 | 0.763 | **0.628** | 0.603 | 0.653 | 0.667 | 0.706 | 0.627 |
| w/o P4 | **0.913** | 0.923 | 0.903 | 0.770 | 0.851 | 0.689 | 0.612 | 0.523 | 0.700 | **0.690** | 0.725 | 0.654 |
| w/o P5 | 0.865 | 0.848 | 0.882 | 0.786 | 0.856 | 0.715 | 0.623 | 0.580 | 0.666 | 0.665 | 0.683 | 0.647 |
| w/o P6 | 0.875 | 0.874 | 0.875 | 0.765 | 0.826 | 0.704 | 0.587 | 0.515 | 0.658 | 0.677 | 0.720 | 0.633 |
| **Maj.** | 0.925 | 0.949 | 0.901 | **0.878** | **0.895** | **0.861** | 0.634 | 0.653 | 0.615 | **0.697** | **0.769** | **0.625** |
| **Fusion** | **0.934** | **0.955** | **0.912** | 0.876 | 0.893 | 0.859 | **0.639** | **0.660** | **0.617** | 0.696 | 0.767 | **0.625** |

| | AEROBLADE | | | CNNSpot | | | CommunityForensics | | | ObjectFormer | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | All | Real | Generated | All | Real | Generated | All | Real | Generated | All | Real | Generated |
| - | 0.852 | 0.834 | 0.870 | **0.918** | **0.978** | **0.857** | 0.861 | 0.850 | 0.872 | 0.903 | 0.910 | 0.815 |

| | DMimageDetection | | | NPR | | | Human (Average) | | | Human (Best) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Public | Real | Generated | All | Real | Generated | All | Real | Generated | All | Real | Generated |
| - | 0.871 | 0.924 | 0.870 | 0.787 | 0.915 | 0.660 | 0.819 | 0.848 | 0.790 | **0.863** | **0.910** | **0.815** |

Table 1: Quantitative Comparison Results on Accuracy. The fusion of P1-6 is giving the best result across all MLLMs, with GPT-4o surpassing traditional methods and the best human on the task. The Maj. row shows the majority vote result.

| Time (s.) | GPT-4o | GPT-4o-mini | Llama 3.2 VI | LLaVA-CoT |
|---|---|---|---|---|
| P0 | 1.4 | 0.9 | 3.2 | 3.4 |
| P1-6 | 47.9 | 25.0 | 67.7 | 69.1 |
| P1-6 + F. | **52.2** | **29.6** | **78.1** | **81.5** |
| P1-6 (Parellel) | 11.4 | 7.2 | 32.8 | 34.0 |
| P1-6 (Parellel) + F. | 16.2 | 6.8 | 37.9 | 41.3 |

Table 2: Evaluation time for different models and different prompts, in seconds per image. The time reported includes image processing.

models relying on OpenAI API calls [50, 51]. Results indicate that P5 and P6 require more time, as they involve larger, multi-turn payloads, but all models maintain reasonable inference times with parallelism, given that natural language output remains rate-limited by token throughput of MLLMs.

## 4.2 AI-Generated Image Detection Accuracy

The results are summarized in Table 1, reporting accuracy for MLLMs: GPT-4o, GPT-4o-mini, Llama 3.2 VI, LLaVa-CoT; traditional methods: AEROBLADE [56], CNNSpot [71], Community-Forensics [53], DMimageDetection [10] and NPR [48]; and human performance. For MLLMs, we report results for baseline (P0), all proposed prompts, majority vote of P1-P6 (Maj.), and fused results. Humans performed well overall (81.9% accuracy), particularly on real images (84.8% accuracy). The best volunteer achieved 4.4% higher accuracy than the average human performance and slightly outperformed the best MLLM baseline (86.3% vs. 85.2% for GPT-4o P0). Other MLLMs performed worse than GPT-4o. While CNNSpot (91.8%) surpassed GPT-4o (85.2%), it operates as a black-box model.

Not all P1–P6 outperform P0, which is expected since these prompts prioritize both accuracy and interoperability. This is a trade-off rather than a direct improvement in both aspects. As a scaled-down version of GPT-4o, GPT-4o-mini performs significantly worse across prompts. They have different trends with different prompts. GPT-4o excels in P3 (reasoning) and P6 (stereotype matching). GPT-4o-mini struggles the most on P3, while in P1 (defect
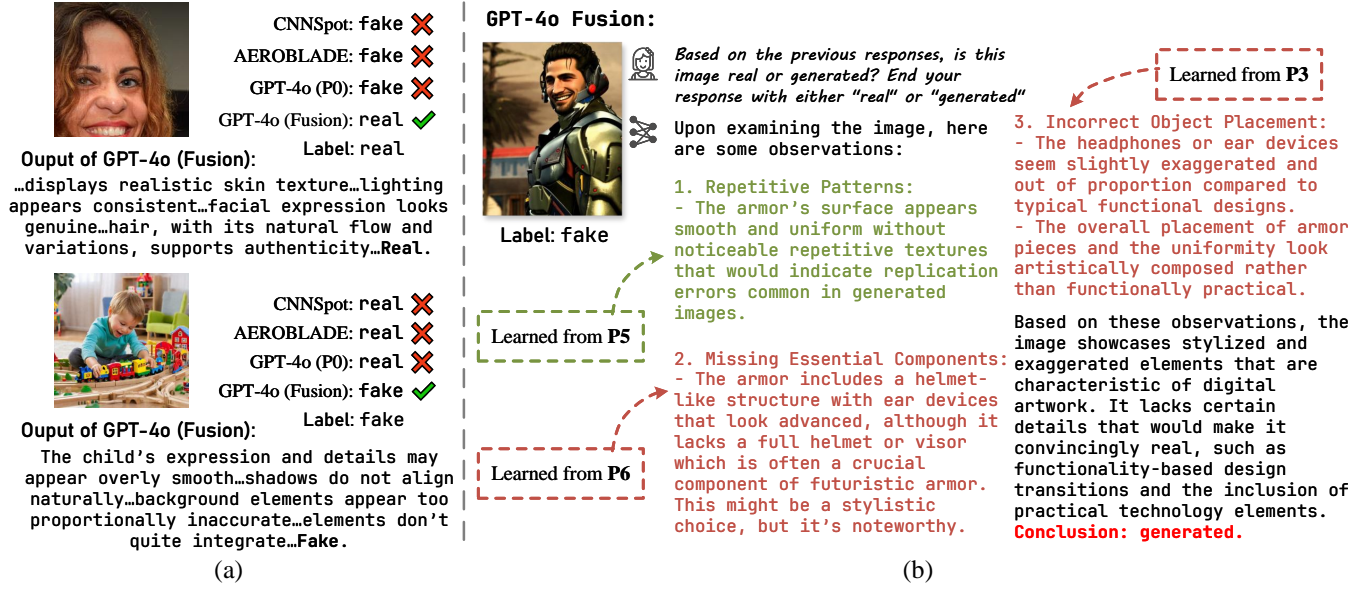
**CNNSpot:** fake ❌
**AEROBLADE:** fake ❌
**GPT-4o (P0):** fake ❌
**GPT-4o (Fusion):** real ✅
**Label:** real

**Ouput of GPT-4o (Fusion):**
…displays realistic skin texture…lighting appears consistent…facial expression looks genuine…hair, with its natural flow and variations, supports authenticity…**Real.**

**CNNSpot:** real ❌
**AEROBLADE:** real ❌
**GPT-4o (P0):** real ❌
**GPT-4o (Fusion):** fake ✅
**Label:** fake

**Ouput of GPT-4o (Fusion):**
The child's expression and details may appear overly smooth…shadows do not align naturally…background elements appear too proportionally inaccurate…elements don't quite integrate…**Fake.**

(a)

**GPT-4o Fusion:**

**Label:** fake

*Based on the previous responses, is this image real or generated? End your response with either "real" or "generated"*

Upon examining the image, here are some observations:

1. Repetitive Patterns:
- The armor's surface appears smooth and uniform without noticeable repetitive textures that would indicate replication errors common in generated images.

Learned from **P5**

2. Missing Essential Components:
- The armor includes a helmet-like structure with ear devices that look advanced, although it lacks a full helmet or visor which is often a crucial component of futuristic armor. This might be a stylistic choice, but it's noteworthy.

Learned from **P6**

Learned from **P3**

3. Incorrect Object Placement:
- The headphones or ear devices seem slightly exaggerated and out of proportion compared to typical functional designs.
- The overall placement of armor pieces and the uniformity look artistically composed rather than functionally practical.

Based on these observations, the image showcases stylized and exaggerated elements that are characteristic of digital artwork. It lacks certain details that would make it convincingly real, such as functionality-based design transitions and the inclusion of practical technology elements.
**Conclusion: generated.**

(b)

**Figure 3: (a) Examples where GPT-4o fusion gives correct result, while CNNSpot, AEROBLADE, and vanilla GPT-4o (P0) fail. (b) A fusion example. GPT-4o can combine responses from P1-6 effectively, drawing conclusions from *reasons* instead of verdicts.**

detection), GPT-4o outperforms GPT-4o-mini, particularly in detecting generated images.

A key finding is the performance gain from response fusion. When using majority vote from verdicts provided by P1 P6, we consistently achieve better accuracy over independent prompts. To fully utilize previous responses, when fusing together for a final verdict, GPT-4o and Llama 3.2 Vision Instruct show an increase in accuracy. While the other two models did not benefit from this process, the aggregated reasons provided by the fusion process are more sensible than individual outputs. When fused, GPT-4o can reach 93.4% accuracy, surpassing CNNSpot. Both majority voting and fusion give more accurate results when compared to P0: *Fused compared to P0, GPT-4o +8.2, GPT-4o-mini + 10.7, Llama 3.2 VI + 7.8, LLaVA-CoT + 10.4.*

These results align with prior research, demonstrating that model ensembling enhances generalization, accuracy, and reduces false positives and false negatives in MLLM-based AI-generated image detection.

## 4.3 Prompt Effectiveness

To assess the effectiveness of our six prompts, we calculated the probability of divergent verdicts across the prompts for a given image, with results presented in Figure 4. Among all real and generated samples, 22.31% received at least one contradictory verdict across prompts from GPT-4o, and 31.44% from LLaVA-CoT. This demonstrates that different prompts can guide the model to focus on different aspects of the image, leading to independent, different decisions.

We also conducted an ablation study by systematically excluding each prompt from the fusion process while preserving the reasoning structure. By analyzing the resulting accuracy, we assess the



**Figure 4: Percentage of cases where prompts give different verdicts. The data is aggregated from all four models evaluated.**

significance of each prompt in enhancing the model's decision-making process. If a prompt is redundant or ineffective, its removal should have minimal impact on the fused accuracy. Conversely, if a prompt is critical, omitting it should lead to a noticeable accuracy drop. For each iteration, we remove one prompt at a time and re-run the fusion queries, comparing the results against the original fused accuracy reported in Table 1. Different models exhibited different sensitivities to prompt removal. For example, GPT-4o and GPT-4o-mini showed the least accuracy drop when omitting P1

| Prompt | GPT-4o | | GPT-4o-mini | | Llama 3.2 VI | | LLaVA-CoT | |
|---|---|---|---|---|---|---|---|---|
| | Accuracy | Rejections ($\downarrow$) | Accuracy | Rejections | Accuracy | Rejections | Accuracy | Rejections |
| P0 ('Generated') | 0.852 | 26 | 0.769 | 23 | 0.561 | 6 | 0.592 | 5 |
| P0 ('Fake') | 0.849 | 45 | 0.761 | 49 | 0.543 | 8 | 0.588 | 5 |

**Table 3: Changing 'fake' to 'generated' in the prompt can lower rejection rates and improve accuracy for all models.**

(defect detection), whereas open-source models retained more accuracy when ablating P1. No ablation configuration outperformed the full fusion model (P1–P6), demonstrating the effectiveness of all proposed prompts.

### 4.4 Qualitative Results on Reasoning

Figure 3 (a) presents two cases where traditional methods and P0 fail, but fused GPT-4o correctly classifies the images by integrating regional observations and general cues. In the first case, GPT-4o with fused prompts identifies realistic skin texture, consistent lighting, genuine facial expression, and naturally flowing hair, leading to a real classification. In the second case, GPT-4o with fused prompts detects overly smooth facial details, unnatural shadows, disproportionate background elements, and poor integration, resulting in a fake classification. These cases demonstrate GPT-4o's holistic reasoning, effectively combining multiple perspectives for accurate detection.

A key advantage of our pipeline is the MLLM fusion stage, which corrects errors from individual prompts (P1–P6). As shown in Figure 3 (b), while P5 incorrectly classifies an image as real, the model re-evaluates observations from P3 and P6, concluding that the image is AI-generated. By integrating previous responses, the model enhances explainability, reinforcing final classifications with well-structured reasoning.

### 4.5 Discussions on MLLM Rejections

MLLMs may generate non-parsable responses or refuse to comment on the given image, termed "rejections," when processing certain images. The rejection rate varies across different models, and the types of images rejected differ accordingly. A frequent justification for these rejections by MLLMs is their inability to provide commentary on images, particularly those depicting faces, due to ethical considerations. These rejections can be mitigated by modifying the prompt, leading to lower rejection rates. Table 3 shows that among the 2,000 images evaluated, replacing the word "fake" with "generated" resulted in a 0.95% (19 samples) reduction in rejections for GPT-4o and 1.3% for GPT-4o-mini. In comparison, open-source models generally exhibit lower rejection rates. Open-source models such as Llama-3.2-Vision-Instruct and LLaVA-CoT have a much lower rejection rate than GPT-4o and -4o-mini.

### 5 Limitations

AI-generated content introduces significant ethical risks, particularly in misinformation, identity manipulation, and fraud. Robust detection systems are essential to address these concerns. Ethically, ensuring transparency and accountability in detection models is

critical, especially in sensitive areas like forensics and law enforcement. These models must be interpretable and free from biases that could lead to unfair outcomes. Privacy and security concerns also arise from the use of AI in surveillance and public discourse, requiring strict safeguards.

While MLLMs show promise in detecting AI-generated images, challenges remain in interpretability and alignment with human perception. Further research is needed to quantify this alignment and assess how human evaluators interpret MLLM-generated explanations in forensic applications. Future research could explore adaptive fusion mechanisms that dynamically adjust reasoning strategies based on image characteristics, as well as adversarial training on high-quality AI-generated images to enhance model robustness.

### 6 Conclusion

In this work, we explore the use of MLLMs for AI-generated image detection, offering a human-interpretable classification method. We introduce six detection paradigms tailored for MLLMs and validate their effectiveness across diverse image generator architectures. A majority vote from these paradigms can surpass traditional classification methods and the most accurate humans. The fusion of model responses, when passed again to the MLLM, allows for aggregated and accurate reasoning of why an image is considered real or AI-generated. Combining diverse heuristics and reasoning improves both accuracy and reliability. Our ablation study confirms that leveraging multiple detection paradigms outperforms single-method classification.

### Acknowledgments

# References

[1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond. *arXiv preprint arXiv:2308.12966* (2023).

[2] Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırlar. 2023. Introducing our Multimodal Models. https://www.adept.ai/blog/fuyu-8b

[3] Xiuli Bi, Bo Liu, Fan Yang, Bin Xiao, Weisheng Li, Gao Huang, and Pamela C. Cosman. 2023. Detecting Generated Images by Real Images Only. arXiv:2311.00962 [cs.CV] https://arxiv.org/abs/2311.00962

[4] Andrew Brock, Jeff Donahue, and Karen Simonyan. 2018. Large Scale GAN Training for High Fidelity Natural Image Synthesis. In *International Conference on Learning Representations*.

[5] Lucy Chai, David Bau, Ser-Nam Lim, and Phillip Isola. 2020. What Makes Fake Images Detectable? Understanding Properties that Generalize. In *Computer Vision – ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXVI* (Glasgow, United Kingdom). Springer-Verlag, Berlin, Heidelberg, 103–120. doi:10.1007/978-3-030-58574-7_7

[6] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C. Lawrence Zitnick. 2015. Microsoft COCO Captions: Data Collection and Evaluation Server. *ArXiv* abs/1504.00325 (2015). https://api.semanticscholar.org/CorpusID:2210455

[7] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 24185–24198.

[8] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. 2018. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 8789–8797.

[9] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. 2020. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8188–8197.

[10] 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 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 1–5. doi:10.1109/ICASSP49357.2023.10095167

[11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, K. Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. *2009 IEEE Conference on Computer Vision and Pattern Recognition* (2009), 248–255. https://api.semanticscholar.org/CorpusID:57246310

[12] Prafulla Dhariwal and Alex Nichol. 2021. Diffusion Models Beat GANs on Image Synthesis. *ArXiv* abs/2105.05233 (2021). https://api.semanticscholar.org/CorpusID:234357997

[13] Prafulla Dhariwal and Alexander Nichol. 2021. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems* 34 (2021), 8780–8794.

[14] Tarik Dzanic, Karan Shah, and Freddie D. Witherden. 2020. Fourier spectrum discrepancies in deep network generated images. In *Proceedings of the 34th International Conference on Neural Information Processing Systems* (Vancouver, BC, Canada) *(NIPS '20)*. Curran Associates Inc., Red Hook, NY, USA, Article 254, 11 pages.

[15] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. 2024. Scaling Rectified Flow Transformers for High-Resolution Image Synthesis. arXiv:2403.03206 [cs.CV] https://arxiv.org/abs/2403.03206

[16] Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 12873–12883.

[17] Joel Frank, Thorsten Eisenhofer, Lea Schönherr, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. 2020. Leveraging frequency analysis for deep fake image recognition. In *Proceedings of the 37th International Conference on Machine Learning (ICML'20)*. JMLR.org, Article 304, 12 pages.

[18] Yueying Gao, Dongliang Chang, Bingyao Yu, Haotian Qin, Lei Chen, Kongming Liang, and Zhanyu Ma. 2025. FakeReasoning: Towards Generalizable Forgery Detection and Reasoning. *arXiv preprint arXiv:2503.21210* (2025). https://arxiv.org/abs/2503.21210

[19] Yifei Gao, Jiaqi Wang, Zhiyu Lin, and Jitao Sang. 2024. AIGCs Confuse AI Too: Investigating and Explaining Synthetic Image-induced Hallucinations in Large Vision-Language Models. arXiv:2403.08542 [cs.CV] https://arxiv.org/abs/2403.08542

[20] Gemini Team. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv:2403.05530 [cs.CL] https://arxiv.org/abs/2403.05530

[21] Gemini Team. 2024. Gemini: A Family of Highly Capable Multimodal Models. arXiv:2312.11805 [cs.CL] https://arxiv.org/abs/2312.11805

[22] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Y. Bengio. 2014. Generative Adversarial Networks. *Advances in Neural Information Processing Systems* 3 (06 2014). doi:10.1145/3422622

[23] Patrick Grommelt, Louis Weiss, Franz-Josef Pfreundt, and Janis Keuper. 2024. Fake or JPEG? Revealing Common Biases in Generated Image Detection Datasets. *ArXiv* abs/2403.17608 (2024). https://api.semanticscholar.org/CorpusID:268691905

[24] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. 2022. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10696–10706.

[25] Zhaojun Guo, Jinghui Lu, Xuejing Liu, Rui Zhao, ZhenXing Qian, and Fei Tan. 2024. What Makes Good Few-shot Examples for Vision-Language Models? arXiv:2405.13532 [cs.CV] https://arxiv.org/abs/2405.13532

[26] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 16000–16009.

[27] Jonathan Ho, Ajay Jain, and P. Abbeel. 2020. Denoising Diffusion Probabilistic Models. *ArXiv* abs/2006.11239 (2020). https://api.semanticscholar.org/CorpusID:219955663

[28] Oleksii Holub. 2022. *Midjourney*. https://www.midjourney.com/home/

[29] Yan Hong, Jianming Feng, Haoxing Chen, Jun Lan, Huijia Zhu, Weiqiang Wang, and Jianfu Zhang. 2025. Wildfake: A large-scale and hierarchical dataset for ai-generated images detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 39. 3500–3508.

[30] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685* (2021).

[31] Ryo Kamoi, Yusen Zhang, Sarkar Snigdha Sarathi Das, Ranran Haoran Zhang, and Rui Zhang. 2024. VisOnlyQA: Large Vision Language Models Still Struggle with Visual Perception of Geometric Information.

[32] Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung Park. 2023. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10124–10134.

[33] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. 2017. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196* (2017).

[34] Tero Karras, Samuli Laine, and Timo Aila. 2018. A Style-Based Generator Architecture for Generative Adversarial Networks. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2018), 4396–4405. https://api.semanticscholar.org/CorpusID:54482423

[35] Amelia Katirai, Noa García, Kazuki Ide, Yuta Nakashima, and Atsuo Kishimoto. 2023. Situating the social issues of image generation models in the model life cycle: a sociotechnical approach. *ArXiv* abs/2311.18345 (2023). https://api.semanticscholar.org/CorpusID:265506098

[36] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.

[37] Tianhong Li, Huiwen Chang, Shlok Kumar Mishra, Han Zhang, Dina Katabi, and Dilip Krishnan. 2023. MAGE: MAsked Generative Encoder to Unify Representation Learning and Image Synthesis. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2142–2152. doi:10.1109/CVPR52729.2023.00213

[38] Yixuan Li, Xuelin Liu, Xiaoyang Wang, Shiqi Wang, and Weisi Lin. 2024. FakeBench: Uncover the Achilles' Heels of Fake Images with Large Multimodal Models. *ArXiv* abs/2404.13306 (2024). https://api.semanticscholar.org/CorpusID:269293612

[39] Yuan-Hong Liao, Rafid Mahmood, Sanja Fidler, and David Acuna. 2024. Reasoning Paths with Reference Objects Elicit Quantitative Spatial Reasoning in Large Vision-Language Models. arXiv:2409.09788 [cs.CV] https://arxiv.org/abs/2409.09788

[40] Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mohammad Shoeybi, and Song Han. 2024. VILA: On Pre-training for Visual Language Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 26689–26699.

[41] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023. Improved Baselines with Visual Instruction Tuning.

[42] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. LLaVA-NeXT: Improved reasoning, OCR, and world knowledge. https://llava-vl.github.io/blog/2024-01-30-llava-next/

[43] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning.

[44] Zhengzhe Liu, Xiaojuan Qi, and Philip H.S. Torr. 2020. Global Texture Enhancement for Fake Face Detection in the Wild. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 8057–8066. doi:10.1109/CVPR42600.2020.00808

[45] Peter Lorenz, Ricard Durall, and Janis Keuper. 2023. Detecting Images Generated by Deep Diffusion Models using their Local Intrinsic Dimensionality. *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)* (2023), 448–459. https://api.semanticscholar.org/CorpusID:259342331

[46] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024. MathVista: Evaluating Mathematical Reasoning of Foundation Models in Visual Contexts. In *International Conference on Learning Representations (ICLR)*.

[47] Shiyin Lu, Yang Li, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, and Han-Jia Ye. 2024. Ovis: Structural Embedding Alignment for Multimodal Large Language Model. *arXiv:2405.20797* (2024).

[48] Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. 2023. Towards Universal Fake Image Detectors that Generalize Across Generative Models. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2023), 24480–24489. https://api.semanticscholar.org/CorpusID:257038440

[49] OpenAI. 2023. *DALL·E 3 System Card*. https://cdn.openai.com/papers/DALL_E_3_System_Card.pdf

[50] OpenAI. 2024. GPT-4o mini: advancing cost-efficient intelligence. https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/.

[51] OpenAI. 2024. GPT-4o System Card. arXiv:2410.21276 [cs.CL] https://arxiv.org/abs/2410.21276

[52] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. 2024. DINOv2: Learning Robust Visual Features without Supervision. arXiv:2304.07193 [cs.CV] https://arxiv.org/abs/2304.07193

[53] Jeongsoo Park and Andrew Owens. 2024. Community Forensics: Using Thousands of Generators to Train Fake Image Detectors. arXiv:2411.04125 [cs.CV] https://arxiv.org/abs/2411.04125

[54] Qwen Team. 2024. QVQ: To See the World with Wisdom. https://qwenlm.github.io/blog/qvq-72b-preview/

[55] Qwen Team. 2025. Qwen2.5-VL. https://qwenlm.github.io/blog/qwen2.5-vl/

[56] Jonas Ricker, Denis Lukovnikov, and Asja Fischer. 2024. AEROBLADE: Training-Free Detection of Latent Diffusion Images Using Autoencoder Reconstruction Error. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2024), 9130–9140. https://api.semanticscholar.org/CorpusID:267335007

[57] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis With Latent Diffusion Models. 10684-10695 pages.

[58] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* 35 (2022), 36479–36494.

[59] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems* 35 (2022), 25278–25294.

[60] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*. 618–626. doi:10.1109/ICCV.2017.74

[61] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. *CoRR* abs/1312.6034 (2013). https://api.semanticscholar.org/CorpusID:1450294

[62] Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020. Denoising Diffusion Implicit Models. *ArXiv* abs/2010.02502 (2020). https://api.semanticscholar.org/CorpusID:222140788

[63] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70* (Sydney, NSW, Australia) *(ICML'17)*. JMLR.org, 3319–3328.

[64] Mingxing Tan and Quoc V. Le. 2019. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. *ArXiv* abs/1905.11946 (2019). https://api.semanticscholar.org/CorpusID:167217261

[65] Ming Tao, Bing-Kun Bao, Hao Tang, and Changsheng Xu. 2023. GALIP: Generative Adversarial CLIPs for Text-to-Image Synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14214–14223.

[66] Ming Tao, Hao Tang, Fei Wu, Xiao-Yuan Jing, Bing-Kun Bao, and Changsheng Xu. 2022. Df-gan: A simple and effective baseline for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16515–16525.

[67] The Llama Team. 2024. The Llama 3 Herd of Models. arXiv:2407.21783 [cs.AI] https://arxiv.org/abs/2407.21783

[68] Aaron Van Den Oord, Oriol Vinyals, et al. 2017. Neural discrete representation learning. *Advances in neural information processing systems* 30 (2017).

[69] Junke Wang, Zuxuan Wu, Jingjing Chen, Xintong Han, Abhinav Shrivastava, Ser-Nam Lim, and Yu-Gang Jiang. 2022. Objectformer for image manipulation detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

[70] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution. *arXiv preprint arXiv:2409.12191* (2024).

[71] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. 2020. CNN-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8695–8704.

[72] Junda Wu, Zhehao Zhang, Yu Xia, Xintong Li, Zhaoyang Xia, Aaron Chang, Tong Yu, Sungchul Kim, Ryan A. Rossi, Ruiyi Zhang, Subrata Mitra, Dimitris N. Metaxas, Lina Yao, Jingbo Shang, and Julian McAuley. 2024. Visual Prompting in Multimodal Large Language Models: A Survey. arXiv:2409.15310 [cs.LG] https://arxiv.org/abs/2409.15310

[73] Penghao Wu and Saining Xie. 2023. V*: Guided Visual Search as a Core Mechanism in Multimodal LLMs. arXiv:2312.14135 [cs.CV] https://arxiv.org/abs/2312.14135

[74] Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. 2024. LLaVA-CoT: Let Vision Language Models Reason Step-by-Step. *ArXiv* abs/2411.10440 (2024). https://api.semanticscholar.org/CorpusID:274116688

[75] Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, and Fei Huang. 2024. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13040–13051.

[76] Shin-Ying Yeh, Yu-Guan Hsieh, Zhidong Gao, Bernard BW Yang, Giyeong Oh, and Yanmin Gong. 2023. Navigating Text-To-Image Customization: From LyCORIS Fine-Tuning to Model Evaluation. *arXiv preprint arXiv:2309.14859* (2023).

[77] Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin, Shuo Liu, Jiayi Lei, Quanfeng Lu, Runjian Chen, Peng Xu, Renrui Zhang, Haozhe Zhang, Peng Gao, Yali Wang, Yu Qiao, Ping Luo, Kaipeng Zhang, and Wenqi Shao. 2024. MMT-Bench: A Comprehensive Multimodal Benchmark for Evaluating Large Vision-Language Models Towards Multitask AGI. arXiv:2404.16006 [cs.CV]

[78] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. 2015. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365* (2015).

[79] Jiarui Zhang, Mahyar Khayatkhoei, Prateek Chhikara, and Filip Ilievski. 2024. Towards Perceiving Small Visual Details in Zero-shot Visual Question Answering with Multimodal LLMs. arXiv:2310.16033 [cs.CV] https://arxiv.org/abs/2310.16033

[80] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3836–3847.

[81] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. 2019. Detecting and Simulating Artifacts in GAN Fake Images. In *2019 IEEE International Workshop on Information Forensics and Security (WIFS)*. 1–6. doi:10.1109/WIFS47025.2019.9035107

[82] Yue Zhang, Ben Colman, Ali Shahriyari, and Gaurav Bharaj. 2024. Common Sense Reasoning for Deep Fake Detection. *ArXiv* abs/2402.00126 (2024). https://api.semanticscholar.org/CorpusID:267365533

[83] Zicheng Zhang, Haoning Wu, Chunyi Li, Yingjie Zhou, Wei Sun, Xiongkuo Min, Zijian Chen, Xiaohong Liu, Weisi Lin, and Guangtao Zhai. 2024. A-Bench: Are LMMs Masters at Evaluating AI-generated Images? *ArXiv* abs/2406.03070 (2024). https://api.semanticscholar.org/CorpusID:270258348

[84] Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. 2023. MMICL: Empowering Vision-language Model with Multi-Modal In-Context Learning. *ArXiv* abs/2309.07915 (2023). https://api.semanticscholar.org/CorpusID:261823391
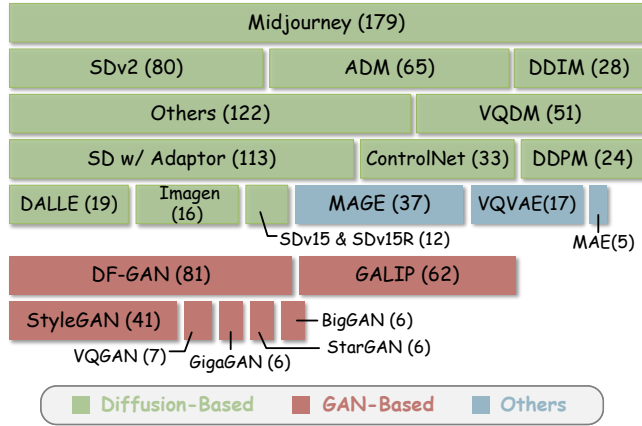
# Appendix

## A  Dataset Distribution



Figure 5: Distribution of AI generation methods among images in the dataset.

To ensure diversity, we curated a dataset containing 2000 images from multiple sources, covering three major architectures and variations of specific models, as detailed below.

**Diffusion-Based Generators.** We sampled from Midjourney [28], SDv2 [57], ADM [13], DDIM [62], DDPM [27], Imagen [58], VQDM [24], ControlNet [80], SD with adaptors (LoRA [30], LyCORIS [76]) and DALLE 3 [49].

**GAN-Based Generators.** We sampled from DF-GAN [66], GALIP [65], StyleGAN [34], VQGAN [16], GigaGAN [32], BigGAN [4] and StarGAN [8, 9].

**Other Generators.** We sampled from VQVAE [68], MAGE [37] and MAE [26].

**Real Images.** We select real images from COCO [6], LAION-5B [59], LSUN Church [78], FFHQ [34], CelebA-HQ [33], ImageNet [11] and AFHQ [9].

## B  Discussions on Model Responses

To accommodate paper length constraints, we present additional results for qualitative analysis in Figure 9. These examples are unfiltered, showcasing the model's typical performance. Our analysis shows that GPT-4o produces varying verdicts across different prompts but maintains consistent reasoning accuracy during the fusion stage. This aligns with our experimental findings, where post-fusion responses achieve the highest accuracy among MLLM-based approaches. Our full prompts can be accessed from this link.

For conciseness, responses are further shortened with the prompt below.

```
<original_response>
---
Summarize the above response into no more than 5 key
    points, where each point should be fewer than 6
    words. Output a numbered list in markdown format.
    Do not include the final verdict in your response.
```

**Key words that model uses.** To qualitatively assess the model's attention tendency, we list the most common key words that are
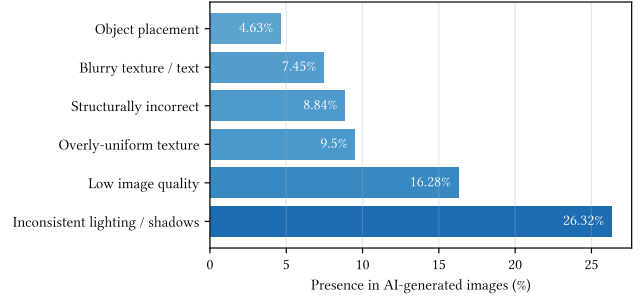


Figure 6: Common key words present in model responses.

seen many times in model responses. While some of the model responses

## C  Experimental Results on Other MLLMs

We evaluate several MLLMs to distinguish AI-generated images from real ones. The results, as shown in Table 4, indicate that Qwen2-VL [70] and QVQ [54] perform relatively well in this task. However, Qwen2-VL exhibits a bias towards classifying images as real, while QVQ, with 72B parameters, can take almost ten minutes for one fusion query on 4x NVIDIA H800 GPU. Qwen2.5-VL-7B-Instruct also gives remarkable results, but it lacks the multi-turn image input chat template, rendering P2 and P4 unusable. Gemini 1.5 Pro [20], as a proprietary model, scored lower than GPT-4o-mini in all three metrics. We selected LLaVA-CoT [74] for its reinforced CoT capabilities and conducted further experiments to explore the potential of open-source MLLMs in AI-generated image detection. Llama 3.2 Vision Instruct, the base model of LLaVA-CoT, is used for comparison.

## D  Discussions on Few-Shot Exemplar Selection

MLLMs exhibit sensitivity to exemplar selection in few-shot learning [25]. To systematically assess its impact, we compare few-shot learning (P3) against the zero-shot paradigm (P0), evaluating how exemplar choice affects model accuracy and decision stability.
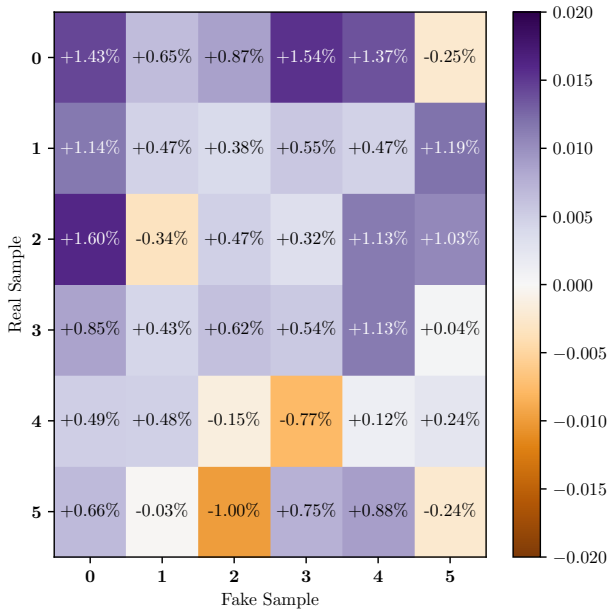
We selected six real and six AI-generated images uniformly. These images were independently annotated by human researchers, identifying artifacts, inconsistencies, and perceptual anomalies that might influence MLLM classification. All 36 possible real-fake combinations were tested on the same dataset as used in Section 4. Additionally, we removed both exemplars from the prompt to create a zero-shot context (P-Zero), serving as the control group.

Figure 8 presents the comparative accuracy of different few-shot sample combinations against the zero-shot baseline. Additionally, Figure 7 provides a grid representation of accuracy variations across different pairings of real and fake exemplars.

From the results, we observe that model accuracy fluctuates depending on the selection of exemplars. The highest observed accuracy gain was +1.60%, while the largest drop was -1.00%. Despite these variations, the overall standard deviation remains moderate, indicating that while sample selection plays a role, GPT-4o exhibits relative stability in performance across different trials.

| Prompt | GPT-4o | | | GPT-4o-mini | | | Gemini 1.5 Pro | | | Kimi-VL-A3B-Instruct | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | All | Real | Generated | All | Real | Generated | All | Real | Generated | All | Real | Generated |
| P0 | 0.852 | 0.855 | 0.849 | 0.769 | 0.910 | 0.628 | 0.758 | 0.906 | 0.610 | 0.585 | 0.622 | 0.548 |
| | MiMo-VL-7B | | | Qwen2-VL-7B-Instruct | | | QVQ-72B-Preview | | | Qwen2.5-VL-7B-Instruct | | |
| | All | Real | Generated | All | Real | Generated | All | Real | Generated | All | Real | Generated |
| P0 | 0.585 | 0.887 | 0.283 | 0.623 | **0.901** | 0.346 | **0.643** | 0.789 | 0.497 | 0.621 | 0.660 | 0.582 |
| | Ovis1.6-Gemma2-9B | | | InternVL2-8B-MPO | | | InternVL2.5-26B | | | MiniCPM-V-2.6 (8B) | | |
| | All | Real | Generated | All | Real | Generated | All | Real | Generated | All | Real | Generated |
| P0 | 0.559 | 0.828 | 0.290 | 0.592 | 0.623 | **0.561** | 0.566 | 0.673 | 0.460 | 0.545 | 0.766 | 0.324 |

Table 4: Experimental result of P0 on different open-source and proprietary models.



Figure 7: Grid representation of accuracy changes for different real-fake exemplar pairs in the few-shot learning setting. Positive values indicate accuracy improvements over the zero-shot baseline (P-Zero), while negative values denote declines.

While sample selection influences performance, the variance remains limited, suggesting GPT-4o exhibits robustness to exemplar variations. In applications demanding high interpretability and reliability, exploring broader sample combinations could further optimize performance.
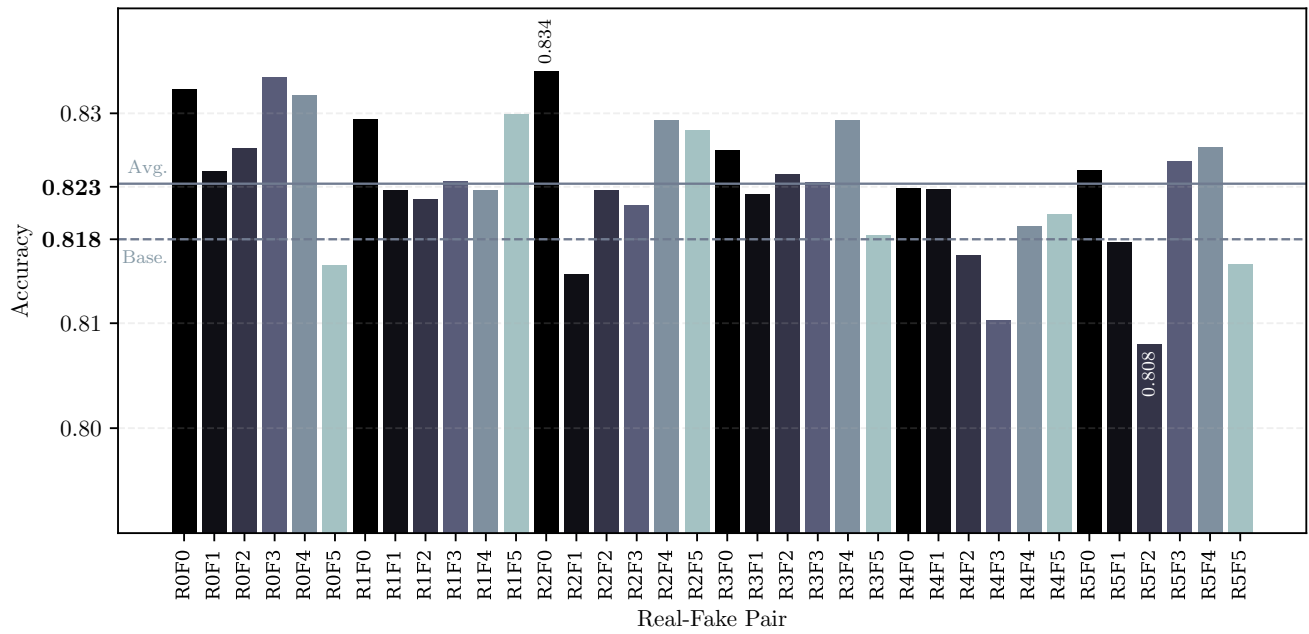
**Figure 8: Impact of real and fake sample selection on the overall accuracy of GPT-4o. The *dashed horizontal line* score represents zero-shot accuracy, while the *solid horizontal line* line denotes the average accuracy across all two-shot pairs.**

| Prompts ▼ | GENERATED | GENERATED | REAL | REAL |
|---|---|---|---|---|
| **P0** Simple / **P1** Defect Query | Generated | Generated | Real | Generated |
| | 1. Consistent logo placement 2. Premium material details 3. Cohesive luxury design 4. Natural hardware elements | 1. Digital texture rendering 2. Uniform lighting patterns 3. Artificial boot appearance 4. Unrealistic scene composition | 1. Text clarity 2. Realistic shadow 3. Texture has natural imperfections | 1. Unnatural surface reflections 2. Artificial edge perfection 3. Abnormal object proportions |
| **P2** ROI | 1. Unnatural strap-to-bag connections 2. Artificial fabric texture rendering 3. Unusual background blur transitions 4. Pattern distortions in branding elements | 1. Too reflective leather texture 2. Boot seem to be floating 3. Background has no texture, only plain color | 1. Natural, non-uniform texture present 2. All essential pouch components visible 3. Natural bokeh and blurring | 1. Window has natural light 2. Too perfect vanity edges 3. Uniformly colored floor 4. Artificial reflections in mirrors |
| **P3** Common Sense | 1. Unnatural strap width 2. Inconsistent logo edges 3. Artificial blur depth 4. Missing hand details | 1. Unusual boot perspective 2. Inconsistent frame depth 3. Unnatural object projection 4. Atypical subject composition | 1. Physical proportions maintain accurate 2. Embroidery distinctly separates from fabric 3. Shape follows expected behavior | 1. Accurate mirror reflections 2. Realistic object proportions 3. Coherent spatial layout |
| **P4** Few-Shot | 1. Too blurry background 2. Colors are unnatural 3. Bag texture connects to clothings 4. Too perfect bag rendering | 1. Natural lighting and shadows 2. Consistent room perspective 3. Realistic furniture textures 4. Coherent wall décor arrangement 5. Good bokeh setup | 1. Natural fabric texture appears 2. Zipper/closure not visible 3. Text follows fabric contours | 1. Professional lighting effects 2. High-end installation precision 3. Consistent spatial relationships |
| **P5** Structural | 1. Natural logo variation 2. Leather material fits the brand 3. Sleek and modern design with clean lines | 1. Over-uniform surface textures 2. Missing shoe detail elements 3. Unnatural display setup | 1. Text has unnatural warping and waviness 2. Inconsistent lighting at text edges | 1. Natural tile patterns 2. Logical bathroom component placement 3. Realistic material textures |
| **P6** Stereotype | 1. Overly smooth pattern details 2. Artificial background blur 3. Missing natural imperfections 4. Inconsistent edge sharpness 5. Uniform texture quality | 1. Mismatched modern setting 2. Artificial boot perfection 3. Atypical display context | 1. Rectangular fabric shape is standard 2. Simple embroidered text appears 3. Plain design, no decorative patterns | 1. Light wall tiling 2. Standard vanity setup 3. Conventional bathroom layout |
| **Fusion** | 1. Pattern distortion at curved edges in the Louis Vuitton monogram pattern 2. Unusually sharp backpack against blurred background 3. Strap connection appears slightly artificial 4. Overly uniform texture on the leather material 5. Unnatural transition between backpack and coat | 1. Unnatural texture perfection 2. Contextual dissonance 3. Perspective irregularities 4. Display setup anomalies | 1. The text appears to be embroidered and follows the natural contours of the fabric. While there is some warping in the text, this is consistent with embroidery on fabric rather than AI generation artifacts. 2. The fabric itself shows natural, non-uniform texture variations - particularly visible in the lighting and shadows. 3. All visible elements are coherent and consistent. | 1. Window light diffusion appears natural, with realistic falloff and shadow casting 2. Cabinet-to-wall transitions show authentic installation details and seams 3. Floor drain and tile layout demonstrate proper construction methodology 4. Material textures (particularly the floor) show subtle, natural variations 5. Hardware mounting and fixture installation display realistic depth relationships |

**Figure 9: More examples with GPT-4o. Responses are summarized. Different colors of the text indicate the verdict for each query (red/green: fake/real).**