



AIGI-Holmes: Towards Explainable and Generalizable AI-Generated Image Detection via Multimodal Large Language Models

Ziyin Zhou^{1*#} Yunpeng Luo^{2*} Yuanchen Wu² Ke Sun¹ Jiayi Ji¹
 Ke Yan² Shouhong Ding² Xiaoshuai Sun^{1†} Yunsheng Wu² Rongrong Ji¹
¹ Key Laboratory of Multimedia Trusted Perception and Efficient Computing,
 Ministry of Education of China, Xiamen University ² Tencent YouTu Lab

{ziyinzhou, skjack}@stu.xmu.edu.cn {xssun, rrji}@xmu.edu.cn jjyxmu@gmail.com

{petterluo, yuanchenwu, kerwinyan, ericshding, simonwu}@tencent.com

<https://github.com/wyczyz/AIGI-Holmes>

Abstract

The rapid development of AI-generated content (AIGC) technology has led to the misuse of highly realistic AI-generated images (AIGI) in spreading misinformation, posing a threat to public information security. Although existing AIGI detection techniques are generally effective, they face two issues: 1) a lack of human-verifiable explanations, and 2) a lack of generalization in the latest generation technology. To address these issues, we introduce a large-scale and comprehensive dataset, Holmes-Set, which includes the Holmes-SFTSet, an instruction-tuning dataset with explanations on whether images are AI-generated, and the Holmes-DPOSet, a human-aligned preference dataset. Our work introduces an efficient data annotation method called the Multi-Expert Jury, enhancing data generation through structured MLLM explanations and quality control via cross-model evaluation, expert defect filtering, and human preference modification. In addition, we propose Holmes Pipeline, a meticulously designed three-stage training framework comprising visual expert pre-training, supervised fine-tuning, and direct preference optimization. Holmes Pipeline adapts multimodal large language models (MLLMs) for AIGI detection while generating human-verifiable and human-aligned explanations, ultimately yielding our model AIGI-Holmes. During the inference stage, we introduce a collaborative decoding strategy that integrates the model perception of the visual expert with the semantic reasoning of MLLMs, further enhancing the generalization capabilities. Extensive experiments on three benchmarks validate the effectiveness of our AIGI-Holmes.

1. Introduction

The rapid evolution of AI technologies like GANs [27, 42] and Diffusion models [21, 71] has made generated images

highly realistic. While beneficial for digital art and film, these technologies also pose risks such as misinformation, privacy breaches, and deepfakes. Recent advancements in diffusion models (e.g., FLUX [4], SD3, SD3.5 [22]) and autoregressive techniques (e.g., VAR [85]) have further complicated detection, highlighting the urgent need for effective AI-generated image detection methods.

Recent studies [45, 61, 83, 89, 91, 98] have demonstrated remarkable advancements in AI-generated image detection. While these improvements are noteworthy, two critical issues limit their application and generalization in real-world scenarios: **1. Lack of explanation:** Current detection models are black boxes (Fig. 1(a)), making their detection results difficult for humans to verify. The lack of human-verifiable explanations leads to unreliable detection results. **2. Lack of generalization:** The rapidly evolving AIGC technologies (Fig. 1(c)) persistently challenge their generalization capabilities. Therefore, developing explainable and generalizable AI-generated image detection algorithms is becoming increasingly urgent.

Recent breakthroughs in Multimodal Large Language Models (MLLMs) [13, 32, 33, 52, 67, 81, 97, 99, 104] offer a promising pathway: their exceptional capabilities in commonsense understanding and natural language generation enable semantic-level analysis of visual content. This makes MLLMs a strong candidate for explainable and generalizable AI-generated image detection algorithms. There have been preliminary explorations of using MLLMs for AI-generated images [5, 43]. However, merely employing MLLMs for binary classification predictions does not fully leverage their potential. We aspire for Multimodal Large Language Models (MLLMs) to emulate the capabilities of the renowned de-

This work was done during an internship at Tencent YouTu Lab.

* Equal Contribution.

†Corresponding Author.

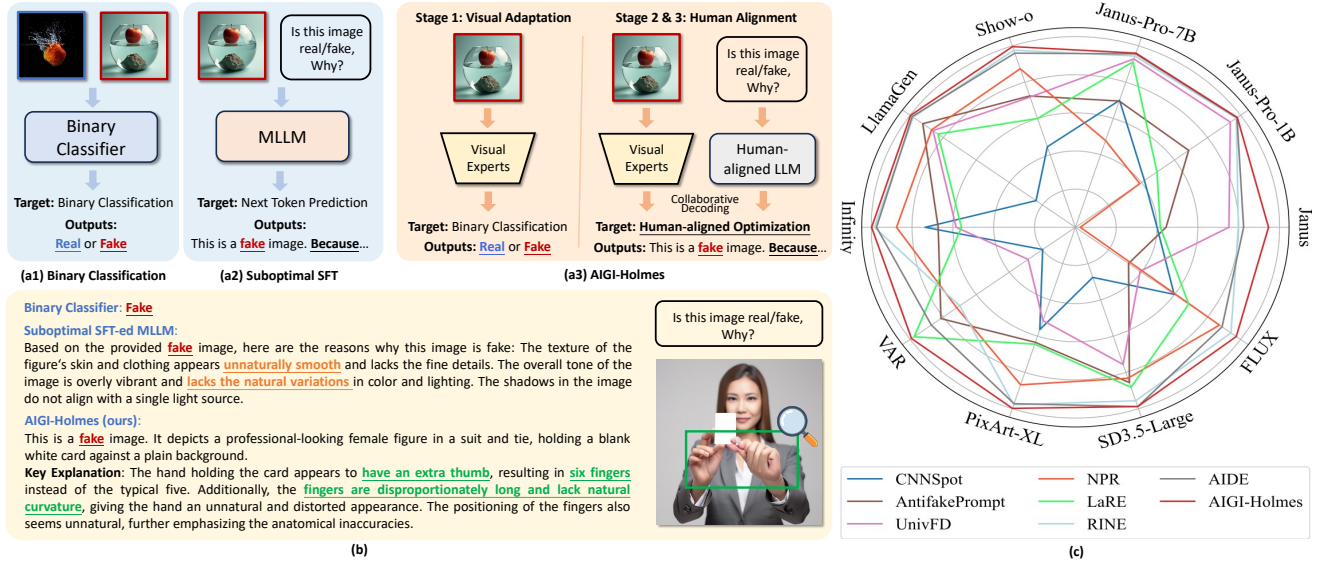


Figure 1. (a): Comparison of AIGI-Holmes with existing methods, (b): A qualitative example to illustrate the effect of AIGI-Holmes, (c): AIGI-Holmes outperforms existing baseline methods on state-of-the-art generators under unseen settings.

tective Sherlock Holmes, not only by accurately identifying the culprit but also by providing precise and corresponding evidence, thereby offering human-verifiable explanations.

However, two key challenges hinder MLLM’s effective use in AI-generated image detection: 1). **Training Data Scarcity:** As shown in Tab. 1, existing AI-generated image detection datasets, such as CNNDetection [88], GenImage [108], and DRCT [7], consist solely of visual modalities and lack instruction fine-tuning datasets suitable for the SFT(Supervised Fine-tuning) phase of MLLMs. Although FakeBench [51] and LOKI [100] make initial attempts in this field, the use of GPT-4o [1] and the high cost of human labor limit the expansion of these datasets, rendering them unsuitable for SFT of MLLMs. 2). **Suboptimal Supervised Fine-Tuning.** Simply training MLLMs on Supervised Fine-Tuning datasets (even with explanations) often leads to sub-optimal performance (as shown in Tab. 3 and Tab. 4). Possible reasons include: 1. MLLMs exhibit somewhat insufficient capabilities in image classification tasks [103] and low-level perception [94], which are often closely related to the generalization of AI-generated image detection. 2. SFT models may mechanically replicate explanation templates without genuinely understanding the underlying causes of artifacts or semantic errors.

To address the issues with the training dataset, we first introduce two key datasets: Holmes-SFTSet and Holmes-DPOSet. The Holmes-SFTSet provides 65K images annotated with explanations across high-level semantic dimensions (e.g., physical inconsistencies, anatomical errors, and text rendering flaws) and low-level artifacts (e.g., overall hue, texture, edges), rigorously refined through cross-model validation and expert-guided filtering to ensure alignment with human-verifiable evidence. To further the critical need

for human-aligned judgment in AI-generated image detection, we introduce the Holmes-DPOSet by constructing contrastive explanation pairs through Positive/Negative prompts on 65K images from the Holmes-SFTSet, along with iterative expert corrections for an additional 4K image explanations, effectively bridging the gap between model perception and human reasoning.

Building upon this training dataset, we propose the Holmes Pipeline, which introduces three key stages to enhance the generalization and interpretability of AIGI detection through systematic MLLM fine-tuning to address the issue of Suboptimal Fine-tuning. The process begins with Visual Expert Pre-training, which leverages the Holmes-SFTSet to rapidly adapt the visual encoder through binary classification, establishing domain-specific feature extraction. Following this, Supervised Fine-Tuning enables MLLMs to not only detect synthetic content but also generate human-verifiable explanations. This stage addresses the “black box” limitation of conventional binary classification approaches. Finally, Human-aligned Direct Preference Optimization [69] utilizes the Holmes-DPOSet. This stage fundamentally reshapes the reasoning patterns of MLLMs by learning from preference samples, ensuring that the interpretative results align with human judgment standards rather than suboptimal fine-tuning. During inference, our collaborative decoding strategy integrates the model perception of the visual expert with the semantic reasoning of MLLMs, creating a dual-channel verification process that enhances the generalizability of our approach. Our method demonstrates superior performance over state-of-the-art approaches in AI-generated image detection, while producing human-aligned explanations that enhance detection reliability.

In summary, our contributions are threefold:

- **Dataset:** We introduce Holmes-SFTSet and Holmes-DPOSet, the first explanation-rich datasets that include human-verifiable evidence through semantic annotations and contrastive preference pairs, addressing the critical training data scarcity in AI-generated image detection.
- **Methodology:** We propose the Holmes Pipeline, a systematic training pipeline for multimodal large language models (MLLMs) that includes visual expert pre-training, explanation-aware supervised fine-tuning, and human-aligned direct preference optimization. This pipeline synergizes model perception with semantic reasoning through novel collaborative decoding during the inference process.
- **Performance:** Our method achieves state-of-the-art detection accuracy on three benchmark datasets and provides human-verifiable explanations, demonstrating superior generalizability and alignment with human judgment.

2. Related Work

2.1. Detection of AI-Generated Fake Images

With the advancement of AI-based image generation technologies, numerous detection methods have emerged, focusing on training on a single AI-generated image method and generalizing to a wide range of AI-generated images. CNNSpot [89] finds that classifiers trained on ProGAN [41] can generalize to unseen GANs using data augmentation. FreDect [24] detects anomalies in the frequency domain of GAN-generated images. Recent methods have explored new perspectives for better generalization. UnivFD [64] proposes to use pre-trained CLIP-ViT [68] features, generalizing to out-of-distribution (OOD) data through nearest neighbor and linear probing. DIRE [91] introduces Diffusion Reconstruction Error (DIRE), distinguishing real images from Diffusion Model (DM)-generated images by measuring reconstruction errors. DRCT [6] introduces Reconstruction Contrastive Learning (RCL), enhancing generalization by generating challenging samples. PatchCraft [107] detects generated images by segmenting them into small patches, applying SRM filters [25], and examining pixel correlations. NPR [83] introduces neighboring pixel relationships, identifying generated content by analyzing local pixel distribution patterns during upsampling. AIDE [98] develops a two-stream framework using both frequency and semantic information. Despite their acceptable performance, these detection methods cannot explain their underlying principles and struggle to generalize against advanced AI-generated techniques.

2.2. Multimodal Large Language Models

Recent advancements in multimodal language large models (MLLMs) have enhanced Image Forgery Detection and DeepFake Detection, enabling interpretable methods. In DeepFake Detection, DD-VQA [104] and FFAA [33] pioneer the use of MLLMs, leveraging human and GPT-

	#Generators	#Image	Explanation	Pref. Data
CNNDetection [88]	11	720K	×	×
GenImage [108]	8	1M+	×	×
DRCT [7]	16	2M	×	×
WildFake [30]	21	3.5M+	×	×
FakeBench [51]	10	6K	✓	×
LOKI [100]	10	3K	✓	×
Holmes-Set	18	65K+4K	✓	✓

Table 1. Comparison between AI-Generated Image Detection Datasets. “Pref. Data” refers to human-aligned preference data.

4o [34] annotated datasets to train models based on InstructBLIP [20] and LLaVA [59]. In Image Forgery Detection, FAKESHIELD [97] combines MLLMs with a visual segmentation model (SAM [44]) for explainable detection (e-IFDL), using a GPT-4o-enhanced dataset (MMTD-Set). ForgerySleuth [81] uses MLLMs with a trace encoder to detect tampering and generate detailed analyses, creating the ForgeryAnalysis dataset. For Text Tampering Detection, TextSleuth [67] builds a large dataset (ETTD) using GPT-4o, employing a two-stage analysis paradigm and fusion mask prompts. These works primarily focus on constructing domain-specific SFT datasets for suboptimal fine-tuning, while neglecting alignment with human preferences. Our method introduces a human-aligned preference dataset, Holmes-DPOSet, and employs Direct Preference Optimization [69] to address the suboptimal fine-tuning issue.

Additionally, some works have explored simple multimodal methods for detecting AI-generated images. Liu *et al.* [57] enhanced the generalization of the detection method by considering the text encoding embeddings of CLIP and introducing frequency-related adapters into the image encoder. AntiFakePrompt [5], Bi-Lora [43], and Jia *et al.* [35] redefine the detection task as a visual question-answering task, combining vision-language models to improve performance on unseen data. Our work not only achieves state-of-the-art generalization detection accuracy but also provides further clues as to why an image is or is not AI-generated.

3. Method

3.1. Data Pipeline

Overview of Data Pipeline. As shown in Fig. 2, our data pipeline consists of four components: Data Source, which includes Data Collection and Image Generation for gathering images to be annotated. These images are then sent to the Automated Annotation section. Subsequently, the annotations are refined in the Preference Modification stage based on expert feedback. Finally, we conduct a Comprehensive Evaluation to assess the model’s generalizability and interpretability. In the following sections, we will introduce the detailed methodologies for each component in sequence.

Data Source. To ensure a diverse dataset with various types of forgeries and defects, as shown in Fig. 2, we initially selected 45K images from existing large-scale AI-generated image detection datasets such as CNNDetection [88], GenIm-

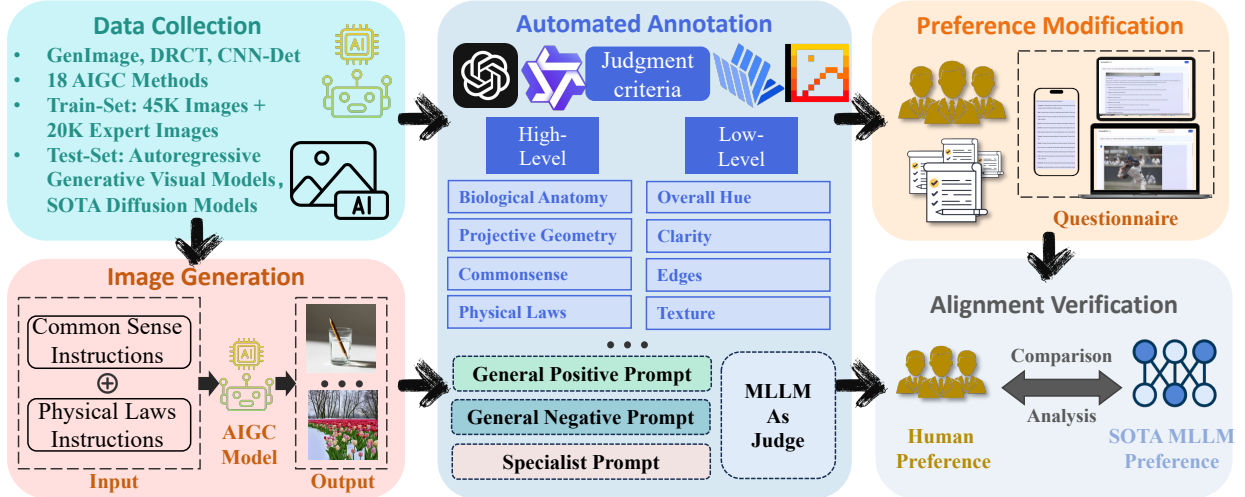


Figure 2. Details of the **Holmes-Set** Construction. The figure illustrates our data pipeline, consisting of four key components: Data Source (including Data Collection and Image Generation), Automated Annotation, Preference Modification (based on human expert feedback), and Comprehensive Evaluation (to assess model generalizability and interpretability).

age [108], and DRCT [7] for image description and forgery explanation. To reduce annotation costs, as in previous work [33, 67, 81, 97], we consider using open-source multimodal large language models (MLLMs) for annotation. A key challenge in constructing annotated explanation datasets by MLLMs for AI-generated images lies in two primary aspects. First, unlike deepfake detection and image forgery detection, AI-generated images often lack corresponding real images, complicating annotation efforts through MLLMs. Second, existing MLLMs exhibit limited capability in analyzing fine-grained forgeries within AI-generated images, as evidenced by recent studies [51]. This limitation likely stems from systematic gaps in domain-specific knowledge related to AI-generated image forensics within the training data of MLLMs. Such knowledge deficiencies may prevent MLLMs from reliably identifying common artifacts in AI-generated content, thereby introducing risks of incomplete or inaccurate annotations.

To mitigate the risk of the inaccurate and incomplete annotation results, inspired by [38, 39], we employed expert-guided methods to filter out 20K images with common AI-generated defects, such as text [10], human bodies [23], human faces [33], projective geometry [72], common sense [26], and physical laws [62], from existing datasets. In simple terms, we use expert small models capable of identifying these defects to filter out the images, and then employ MLLMs to determine the presence of these defects. Additionally, due to the lack of attention to commonsense and physical law defects in existing datasets, we supplement the dataset with images generated from [26, 62], as shown in the Image Generation part of Fig. 2. For these images containing common AI-generated defects, we can design targeted prompts to avoid potential annotation hallucinations, thereby reducing the need for extensive human

labor in the annotation process. *Further details of the image collection process are elaborated in the Appendix C.*

Automated Annotation. We design the automated annotation system named Multi-Expert Jury, comprising four open-source multimodal large language models (MLLM-Experts): Qwen2VL-72B [87], InternVL2-76B [16], InternVL2.5-78B [15], and Pixtral-124B [63], to ensure high-quality data annotation. The system employs three tailored prompts: 1. **General Positive Prompt** as shown in Fig. 6 and Fig. 7: Annotates 45K randomly selected images using high-level (e.g., anatomy, physical laws) and low-level (e.g., texture, clarity) criteria (see Fig. 1). 2. **General Negative Prompt** as shown in Fig. 6 and Fig. 7: Generates adversarial annotations by asking questions contradicting image authenticity, forming natural positive-negative pairs with the General Positive Prompt annotations to construct the Direct Preference Optimization (DPO) dataset \mathcal{D}_1 . 3. **Specialist Prompt** as shown in Fig. 9, Fig. 10, Fig. 11, Fig. 12, Fig. 13: Guides MLLM-Experts to annotate 20K expert-filtered defective images, focusing on generation defects (e.g., commonsense). To ensure annotation quality, we adopt an MLLM-as-a-judge approach [8], where MLLM-Experts cross-evaluate each annotation. Only annotations with top consensus scores are retained in the dataset.

Preference Modification. The constructed SFT dataset can be directly used for coarse-grained alignment. However, there may be instances where models mechanically replicate explanation templates without genuinely understanding the underlying causes of artifacts or semantic errors. To further address this issue and narrow the gap between model perception and human reasoning, we introduce the Holmes-DPO dataset. We manually annotate an additional 2K data samples from the same sources as the training set. During this annotation process, we provide the outputs of the SFT model

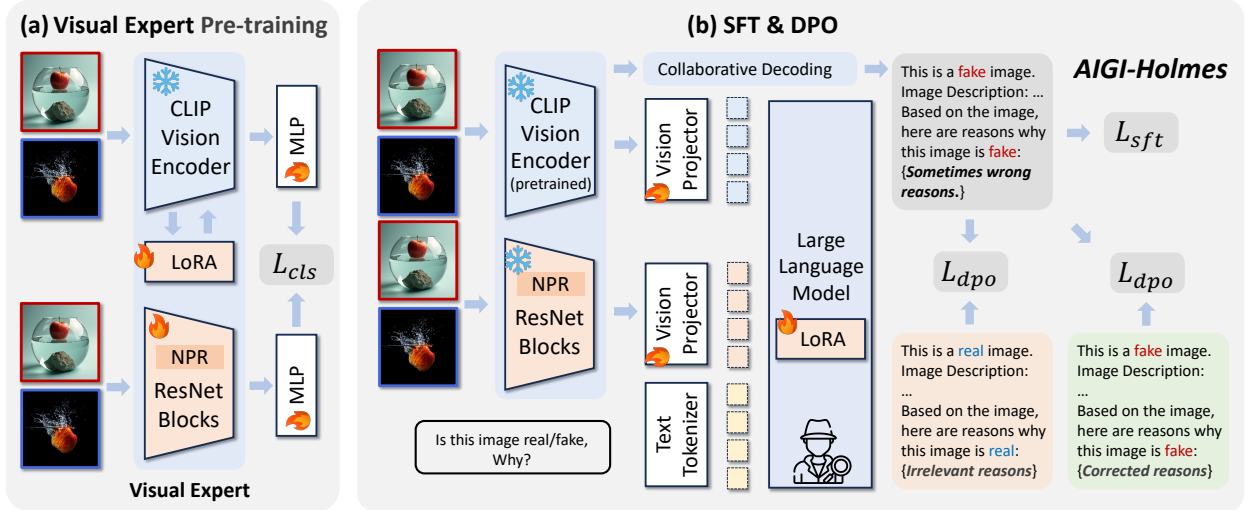


Figure 3. Overview of AIGI-Holmes. We enhance LLaVA [59] with NPR [83] visual expert \mathcal{R} and the Holmes Pipeline, featuring three training stages: Visual Expert Pre-training, SFT, DPO, and a collaborative decoding strategy during inference.

and ask humans to offer modification suggestions, such as supplementary correct information, removal of incorrect or irrelevant explanations, and other modification opinions. We then use Deepseek-V3 [55], an advanced open-source large language model, to modify the original responses of the SFT model based on the suggestions of human experts, resulting in new, more human-aligned correct explanations. The human experts consist of three authors and two annotators who are strictly guided by us (we reference [38, 39] to create annotation documents for training and guiding the annotators). In addition, we use MLLMs and specialist prompts to modify responses to 2K filtered images with known specific defects. We add the pairs of samples before and after modification to the DPO dataset \mathcal{D}_2 .

Comprehensive Evaluation. We evaluate the model’s capabilities from two perspectives: detection and explanation. **1. For detection**, we output the probabilities of the real/fake tokens to calculate Accuracy (Acc.) and Average Precision (A.P.). **2. For explanation**, we use 1K test samples containing Ground Truth, which have been reviewed by annotators for deficiencies in explanations from a professional perspective and corrected uniformly by the Deepseek-V3 model. We calculate metrics such as BLEU [65], CIDEr [86], METEOR [3], and ROUGE [53] to measure the quality of the explanatory text output by the model. Additionally, we employ multimodal large model scoring and human preference evaluation methods for assessment. For MLLM scoring, we refer to [8], using a prompt that considers relevance, accuracy, comprehensiveness, creativity, and granularity to compare and score the model’s responses. For human evaluation scores, we sample 10 images for each type of forgery from \mathcal{P}_3 , resulting in a total of 100 images. The comparative models generate explanations for these test images, and we use pairwise comparison, as referenced in [17], to calculate the ELO ratings for each model’s explanations.

3.2. Overview of AIGI-Holmes

As shown in Fig. 3, for the architecture, we augment the original multimodal method LLaVA [59] with a low-level information NPR [83] visual expert \mathcal{R} . For the training methodology, we introduce the Holmes Pipeline, which includes three key training stages: The visual Expert Pre-training Stage, the SFT Stage, and the DPO Stage. During the final inference, we employ a collaborative decoding strategy to merge the predictions from the visual expert and the language model.

3.3. Architecture

Many existing methods use LLaVA [59], a Multimodal Large Language Model (MLLM) with strong multimodal understanding capabilities, as a baseline. However, this architecture faces several challenges: 1. Limited performance when dealing with classification problems [103]. 2. Inefficiency in handling low-level information [94], which is crucial for AIGI detection. To address these challenges, we augment the CLIP visual encoder \mathcal{F} used in LLaVA with a low-level information NPR [83] visual expert \mathcal{R} . Given an image input X_{img} , we extract visual features f_{img} and f_{npr} through \mathcal{F} and \mathcal{R} , respectively, as shown in the following equations:

$$X_{npr} = \text{NPR}(X_{img}); f_{npr} = \mathcal{R}(X_{npr}); f_{img} = \mathcal{F}(X_{img}). \quad (1)$$

Subsequently, the features are injected into the large language model through a projector and text embedding f_t , producing the final output \mathbf{H} . This process is represented by the following equations:

$$\mathbf{H} = \text{LLM}(\text{proj}([f_{img}, f_{npr}]), f_t). \quad (2)$$

3.4. Holmes Pipeline

Visual Expert Pre-training Stage: To transform the MLLM into an expert in the AIGI detection domain, we need to ensure that the visual expert provides a certain level of generalization capability and detection accuracy. Therefore, before

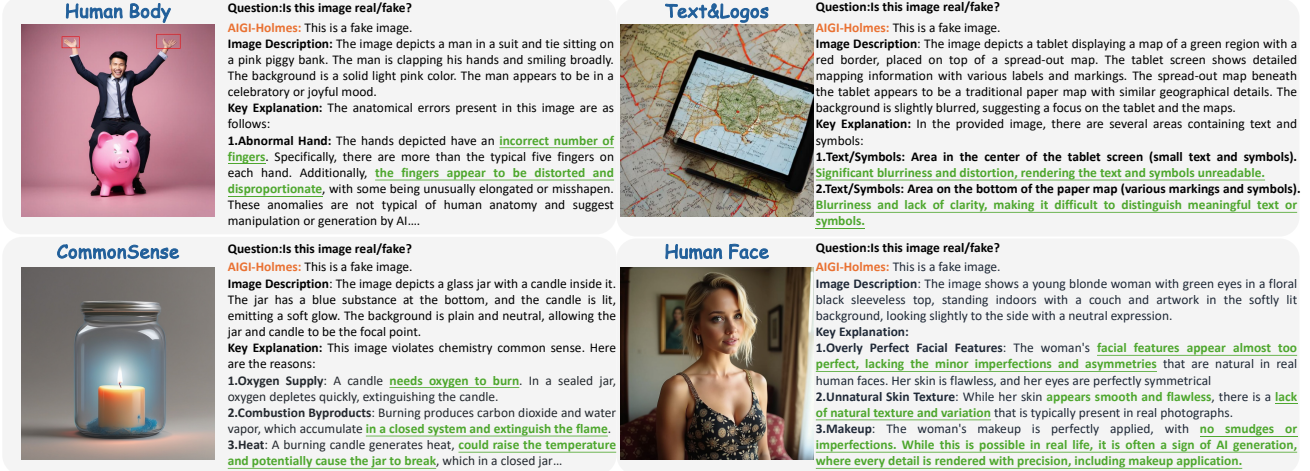


Figure 4. Qualitative results of AIGI-Holmes on AI-Generated images.

the general LLaVA training paradigm, we introduce a pre-training stage for the visual expert. We employ LoRA [31] to efficiently fine-tune CLIP-ViT-L/14, setting $r = 4$ and $\alpha = 8$. The CLS features f_{cls} extracted by CLIP are fed into an MLP to obtain classification results. Simultaneously, we independently adjust the NPR-based ResNet. Following [83], we only use the first two layers of ResNet to extract the features f_{npr} . Similarly, an MLP is used to obtain classification results. For ResNet, we perform full parameter fine-tuning. We use binary cross-entropy loss l_{bce} to adjust the visual expert, as shown in the following equation:

$$\begin{aligned} \mathbf{y}_{clip} &= \text{MLP}(f_{cls}), & \mathbf{y}_{npr} &= \text{MLP}(f_{npr}), \\ l_{clip} &= l_{bce}(\mathbf{y}_{clip}, \mathbf{y}), & l_{npr} &= l_{bce}(\mathbf{y}_{npr}, \mathbf{y}). \end{aligned} \quad (3)$$

SFT Stage: After obtaining the pre-trained visual expert, we integrate it into the large language model. To guide the model to output explanations related to whether an image is AI-generated, we perform SFT on the Holmes-SFT dataset. In this stage, consistent with previous methods, we freeze the visual expert while keeping the linear projector and the LoRA components of the large language model trainable. We optimize these parameters using the autoregressive text loss l_{txt} :

$$l_{txt} = l_{ce}(\mathbf{H}, \mathbf{H}_{txt}). \quad (4)$$

DPO Stage: To further enhance the ability of the SFT model to produce high-quality, human-aligned explanations, we perform human-aligned direct preference optimization on the previously constructed Holmes-DPOSet. Specifically, we sample preference pairs $\{y_w\}, \{y_l\} \sim \mathcal{D}$ from the Holmes-DPOSet, where $\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2$. We then optimize the model using the DPO loss as shown in Eq. 5.

$$\begin{aligned} \mathcal{L}_{\text{DPO}}(\phi) &= -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \left(\log \frac{\pi_{\phi}(y_w|x)}{\pi_{\text{base}}(y_w|x)} - \log \frac{\pi_{\phi}(y_l|x)}{\pi_{\text{base}}(y_l|x)} \right) \right) \right], \end{aligned} \quad (5)$$

where σ and β have the same meanings as in the first round. $\pi_{\phi}(y|x)$ denotes the policy model to be optimized in the second round, with parameters ϕ . $\pi_{\text{base}}(y|x)$ is the model after the first round of preference optimization. \mathcal{D}_2 represents the dataset for the second round of preference optimization. During the preference optimization stage, we keep the visual expert frozen, allow the projector to be trainable, and use LoRA to train the large language model part.

Inference Stage: We propose **Collaborative Decoding**, which aims to utilize both our MLLM and its pre-trained expert to jointly decide the authenticity of an image during inference, thereby enhancing the generalization and detection accuracy of AIGI-Holmes. Specifically, we adjust the logit values of the tokens corresponding to “real” and “fake” in the model’s output. We denote the logit value for “real” as $\text{logit}(\mathbf{y} = 0)$ and the logit value for “fake” as $\text{logit}(\mathbf{y} = 1)$, for $k \in \{0, 1\}$ as shown in Eq. 6:

$$\begin{aligned} \text{logit}_{new}(\mathbf{y} = k) &= \alpha \cdot \text{logit}_{raw}(\mathbf{y} = k) \\ &\quad + \beta \cdot \text{logit}(\mathbf{y}_{clip} = k) + \gamma \cdot \text{logit}(\mathbf{y}_{npr} = k), \end{aligned} \quad (6)$$

where $\alpha = 1$, $\beta = 1$, and $\gamma = 0.2$ are the weights assigned to the three prediction results. Through collaborative decoding involving the pre-trained visual expert, we retain the predictions of the MLLM while preventing it from overfitting to existing forgery types, thereby improving the detection accuracy of the MLLM in unseen domains.

4. Experiment

4.1. Experimental Setup

Datasets. To comprehensively evaluate the generalization capabilities of existing methods, we conducted experiments under three settings: **Protocol-I**, **Protocol-II**, and **Protocol-III**. For **Protocol-I**, we trained on the 4-class (car, cat, chair, horse) subset of the CNNDetection dataset, which was widely used in earlier studies, and evaluated the detector

Method	Janus		Janus-Pro-1B		Janus-Pro-7B		Show-o		LlamaGen		Infinity		VAR		PixArt-XL		SD3.5-Large		FLUX		Mean	
	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.	Acc.	A.P.
CNNSpot	70.0	86.0	70.9	85.8	85.0	93.6	72.2	86.0	61.9	71.4	86.8	94.6	59.9	75.0	78.2	90.1	63.8	81.1	79.9	92.0	72.9	85.6
AntifakePrompt	72.2	87.4	84.3	94.0	84.8	93.1	86.2	95.5	96.2	99.4	83.6	94.1	90.7	95.6	81.7	92.8	92.8	97.8	66.1	80.8	83.9	93.1
UnivFD	87.6	97.8	96.9	99.5	96.4	99.5	85.9	97.4	93.1	98.6	79.2	96.2	64.3	85.9	75.7	94.4	87.8	97.8	69.6	91.4	83.6	95.9
NPR	51.2	55.9	69.5	75.1	73.9	77.9	93.7	99.6	93.5	99.4	93.8	99.9	85.9	91.2	93.4	99.1	91.6	97.7	93.6	99.5	84.0	89.5
LaRE	70.8	99.3	74.7	97.5	95.6	99.7	80.0	99.0	91.6	99.6	77.9	99.6	98.8	100.0	82.2	99.7	94.1	99.5	84.3	99.0	85.0	99.3
RINE	89.9	98.3	98.7	99.9	97.2	99.6	98.8	99.9	99.1	100.0	99.2	99.9	85.0	97.9	98.9	99.8	97.8	99.7	97.1	99.7	96.2	99.5
AIDE	91.2	99.1	98.9	99.9	97.8	99.8	98.0	99.8	99.4	100.0	98.7	99.9	93.6	99.3	98.6	99.9	99.4	100.0	94.4	99.5	97.0	99.7
AIGI-Holmes	97.3	99.9	99.0	99.9	98.0	99.9	99.8	99.9	99.9	100.0	99.9	100.0	99.6	100.0	99.9	100.0	99.4	99.9	98.7	99.7	99.2	99.9

Table 2. Evaluation on the \mathcal{P}_3 . All baseline results are trained on our training set to ensure a fair comparison.

on the general and comprehensive benchmark AIGCDetect-Benchmark [106] (\mathcal{P}_1). For **Protocol-II**, we trained on the training set proposed by AntiFakePrompt [5] and tested on its proposed test set, which includes 18 types of forgeries. For **Protocol-III**, we trained on the dataset containing various Diffusion methods proposed in Sec. 2 and tested on images generated by the latest unseen autoregressive visual generation models Janus [93], Janus-Pro [12], VAR [84], Infinity [29], Show-o [96], LlamaGen [80], and the state-of-the-art diffusion models PixArt-XL [9], FLUX [4], SD3.5 [22] (\mathcal{P}_3). Each of the above test sets contains 5K real images from COCO [54] and 5K generated images corresponding to the generation methods.

Implementation Details. During the pre-training phase of the visual expert, we fine-tune CLIP-ViT/L-14 using LoRA ($r=4$, $\alpha=8$) and fully fine-tune the first two layers of ResNet with NPR as input. The training is conducted for 5 epochs with a batch size of 32. During the SFT phase, we fine-tune LLaVA1.6-mistral-7B [58] using LoRA ($\text{rank}=128$, $\alpha=256$) while fully training the domain label generator. This model is trained for 3 epochs with a learning rate of $5e-5$, a batch size of 16, and a gradient accumulation step of 1. During the DPO phase, we use LoRA ($\text{rank}=48$, $\alpha=96$) with a learning rate of $5e-7$, a batch size of 4, a gradient accumulation step of 2, and $\beta=0.1$, training for 2 epochs.

Comparison Baselines. We compare various baselines in the main text including CNNSpot [89], NPR [83], AntiFakePrompt [5], LaRE [61], RINE [45], AIDE [98]. *Additional baseline methods will be introduced in the Appendix D.1.*

4.2. Comparisons with SOTA Detection Methods

Protocol-III. The quantitative results in Tab. 2 show the classification accuracy of various methods and generators within the range of \mathcal{P}_3 . All methods were retrained on our proposed training set to ensure a fair comparison. The test images were generated by unseen state-of-the-art autoregressive visual models and diffusion models. On this challenging benchmark, AIGI-Holmes achieved SOTA results, with accuracy improvements of 15.2%, 3.0%, and 2.2% over the previous best methods NPR, RINE, and AIDE, respectively. For the detection accuracy on the best autoregressive visual generation techniques and diffusion model representatives VAR and FLUX, our method surpassed the best methods by 6.0% and 1.8%, respectively. These three different training settings emphasize the excellent generalization ability of our proposed AIGI-Holmes. *The results of Protocol-I and*

Protocol-II can be found in the Appendix D.3.

4.3. Comparisons with SOTA MLLMs

As demonstrated in Tab. 3, we conduct a quantitative comparison between the textual explanations generated by AIGI-Holmes and those produced by state-of-the-art (SOTA) multimodal large language models (MLLMs). To ensure a fair comparison, the explanations from the baseline models are obtained under the General Positive Prompt query. Our method achieves the highest metrics across nearly all evaluated aspects. For instance, compared to the state-of-the-art closed-source model GPT-4o, our model’s output achieves a BLEU-1 score of 0.622, which is an improvement of 0.189 over GPT-4o’s score of 0.433. Additionally, for human pairwise scoring, our model achieves an ELO Rating of 11.420, surpassing GPT-4o’s score of 10.271 by 1.149 points. *A selection of model outputs is illustrated in the Appendix E.* This result indicates that existing MLLMs have the potential to provide reasonable explanations for AI-generated images. However, due to the lack of targeted downstream task training and the relatively small number of synthetic images included in the dataset, they are unable to perform more precise analyses in the task of explainable AIGI detection. Our method effectively fills this gap.

4.4. Robustness Evaluation

In real-world scenarios, AI-generated images often encounter unpredictable perturbations during dissemination, which can lead to the failure of existing AI detectors. Within the range of \mathcal{P}_3 , we applied several common perturbations found in real-world scenarios: JPEG compression (QF=75, QF=70), Gaussian blur ($\sigma=1$, $\sigma=2$), and downsampling ($\times 0.5$). As shown in Tab. 5, the performance of all methods significantly declines under these distortions. However, AIGI-Holmes achieves higher detection accuracy compared to other baseline methods in these challenging scenarios. A possible reason is our use of the pre-trained CLIP method, which, as indicated in [18], demonstrates good robustness when used as a backbone for AI image detection methods. Additionally, MLLM focuses more on high-level semantic features, reducing the model’s reliance on low-level artifacts that are crucial for AI-generated image detection in other methods. These artifacts are often susceptible to unpredictable perturbations in real-world scenarios. Additionally, as shown in Fig. 5, under these degradation conditions, the evaluation metrics for model explanations such as BLEU-1,

MLLM	Automatic Metrics				MLLM-as-Judge Evaluation				Human. ELO Ratings
	BLEU-1	ROUGE-L	METEOR	CIDEr	Qwen-72B	InternVL-76B	InternVL-78B	Pixtral-124B	
Qwen2VL-72B	0.314	0.227	0.292	0.003	3.874	3.612	4.002	3.163	8.432
InternVL2-76B	0.362	0.224	0.289	0.006	4.042	3.807	4.006	3.463	10.111
InternVL2.5-78B	0.275	0.221	0.293	0.007	4.012	3.531	3.954	3.101	8.623
Pixtral-124B	0.428	0.270	0.302	0.010	3.967	3.990	4.140	4.213	10.472
GPT-4o	0.433	0.308	0.306	0.005	4.102	4.010	4.032	4.010	10.271
AIGI-Holmes(w/o DPO)	0.445	0.315	0.317	0.023	4.119	3.918	4.150	4.130	10.670
AIGI-Holmes(w/ DPO)	0.622	0.375	0.311	0.107	4.196	4.011	4.189	4.227	11.420

Table 3. A comprehensive comparison of the explanations for AI-generated images between pre-trained SOTA MLLMs and AIGI-Holmes. The abbreviations “w/o” stands for “without”, and “w/” stands for “with”.

VEP-S	DPO	CD	\mathcal{P}_1	\mathcal{P}_3
			83.3	90.1
	✓		84.8	92.3
✓			86.8	97.2
✓	✓		87.4	97.6
✓		✓	90.8	98.9
✓	✓	✓	93.2	99.2

Table 4. Ablation Study of core model components.

Method	JPEG Compression		Gaussian Blur		Resize
	QF=75	QF=70	$\sigma = 1.0$	$\sigma = 2.0$	$\times 0.5$
CNNSpot	63.5	62.4	64.5	61.7	59.9
NPR	52.2	51.6	56.8	53.4	74.3
UnivFD	84.7	84.0	81.0	74.9	86.3
LaRE	62.0	63.0	54.3	54.2	51.1
AntifakePrompt	80.1	79.7	78.2	77.6	74.5
AIDE	92.8	92.3	91.9	90.7	89.2
RINE	92.4	91.1	94.2	92.8	92.3
AIGI-Holmes	99.0	98.7	98.3	97.9	95.9

Table 5. Robustness of Classification Accuracy on JPEG Compression, Gaussian Blur and Resize of AIGI-Holmes. The classification accuracy (%) averaged over 10 test sets in \mathcal{P}_3 with specific perturbation.

ROUGE-L, METEOR, and CIDEr did not exhibit significant declines. This indicates that the explanations generated by the model remain focused on high-level semantic information related to the image content and are not overly affected by these degradation conditions.

4.5. Ablation Study

We conduct ablation experiments on the main innovative methods of AIGI-Holmes: Visual Expert Pre-training Stage (VEP-S), DPO, and Collaborative Decoding (CD), as shown in Tab. 4. The results demonstrate a significant improvement in accuracy when using the Visual Expert Pre-training Stage compared to the original llava paradigm, with detection accuracy increasing by 3.5% and 7.1% on \mathcal{P}_1 and \mathcal{P}_3 , respectively. Adding DPO and Collaborative Decoding on this basis further enhances accuracy. Specifically, DPO helps improve the model’s accuracy by 0.6% and 0.4% on \mathcal{P}_1 and \mathcal{P}_3 , respectively, while Collaborative Decoding in conjunction with the Visual Expert Pre-training Stage boosts accuracy by 4.0% and 1.7%. The combined use of these approaches results in improvements over previous combinations, achieving approximately a 10% increase in accuracy compared to the baseline method. These experiments demonstrate the effectiveness of our design. After employing the DPO stage, a comprehensive improvement in the quantitative results of model output explanations is observed, particularly with an increase of 0.75 points in human-scored ELO Ratings compared to the SFT model. This indicates that DPO is an

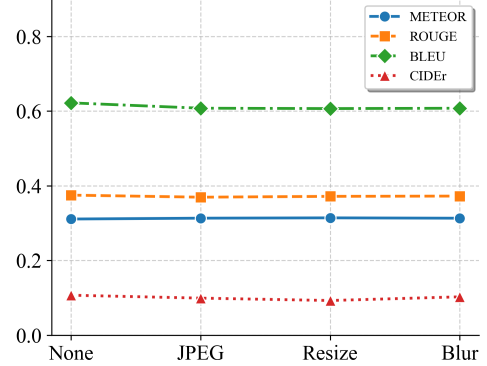


Figure 5. Robustness of the explanation on JPEG Compression (QF=70), Gaussian Blur ($\sigma = 2$), and Resize ($\times 0.5$) of AIGI-Holmes.

effective post-training method for enhancing the quality of model output explanations to better align with human preferences. *Results of other ablation experiments can be found in the Appendix D.4.*

4.6. Qualitative Results

Fig. 4 showcases the detection and explanation results generated by AIGI-Holmes on various AI-generated images. It can be observed that our method produces precise explanations for defects across different AI-Generation modes. *More visual results can be found in the Appendix E.*

5. Conclusion

In this work, we introduce Holmes-Set, the first explanation-rich dataset with human-verifiable semantic annotations and contrastive preference pairs, addressing the critical data scarcity in AI-generated image detection. Besides, the proposed Holmes Pipeline systematically integrates visual expert pre-training, explanation-aware fine-tuning, and human-aligned preference optimization, synergizing model perception with human reasoning. Extensive experiments demonstrate our AIGI-Holmes’ state-of-the-art detection accuracy and human-aligned interpretability. These contributions advance explainable and generalizable AIGI detection for rapidly evolving AIGC scenarios. We hope our dataset and methodology will inspire future research toward building more trustworthy AI-generated image detection systems.

Acknowledgements

This work was supported by the National Science Fund for Distinguished Young Scholars (No.62025603), the National Natural Science Foundation of China (No. U21B2037, No. U22B2051, No. U23A20383, No. U21A20472, No. 62176222, No. 62176223, No. 62176226, No. 62072386, No. 62072387, No. 62072389, No. 62002305 and No. 62272401), and the Natural Science Foundation of Fujian Province of China (No. 2021J06003, No.2022J06001).

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2
- [2] AI@Meta. Llama 3 model card. 2024. 3
- [3] Satandeep Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005. 5
- [4] Black-Forest-Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 1, 7, 2
- [5] You-Ming Chang, Chen Yeh, Wei-Chen Chiu, and Ning Yu. Antifakeprompt: Prompt-tuned vision-language models are fake image detectors. *arXiv preprint arXiv:2310.17419*, 2023. 1, 3, 7, 2, 4
- [6] 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*. 3
- [7] Baoying Chen, Jishen Zeng, Jianquan Yang, and Rui Yang. DRCT: Diffusion reconstruction contrastive training towards universal detection of diffusion generated images. In *Proceedings of the 41st International Conference on Machine Learning*, pages 7621–7639. PMLR, 2024. 2, 3, 4
- [8] Dongping Chen, Ruoxi Chen, Shilin Zhang, Yaochen Wang, Yinyao Liu, Huichi Zhou, Qihui Zhang, Yao Wan, Pan Zhou, and Lichao Sun. Mllm-as-a-judge: Assessing multimodal llm-as-a-judge with vision-language benchmark. In *Forty-first International Conference on Machine Learning*, 2024. 4, 5
- [9] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- σ : Weak-to-strong training of diffusion transformer for 4k text-to-image generation, 2024. 7, 2
- [10] Jingye Chen, Yupan Huang, Tengchao Lv, Lei Cui, Qifeng Chen, and Furu Wei. Textdiffuser: Diffusion models as text painters. *Advances in Neural Information Processing Systems*, 36, 2024. 4
- [11] Shen Chen, Taiping Yao, Hong Liu, Xiaoshuai Sun, Shouhong Ding, Rongrong Ji, et al. Diffusionfake: Enhancing generalization in deepfake detection via guided stable diffusion. *Advances in Neural Information Processing Systems*, 37:101474–101497, 2024. 1
- [12] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Januspro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025. 7, 2
- [13] Yize Chen, Zhiyuan Yan, Siwei Lyu, and Baoyuan Wu. X^2 -dfd: A framework for eXplainable and eXtendable deepfake detection. *arXiv preprint arXiv:2410.06126*, 2024. 1, 5
- [14] Zhongxi Chen, Ke Sun, Ziyin Zhou, Xianming Lin, Xiaoshuai Sun, Liujuan Cao, and Rongrong Ji. Diffusionface: Towards a comprehensive dataset for diffusion-based face forgery analysis. *arXiv preprint arXiv:2403.18471*, 2024. 1
- [15] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 4
- [16] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024. 4
- [17] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. In *Forty-first International Conference on Machine Learning*, 2024. 5, 4
- [18] Davide Cozzolino, Giovanni Poggi, Riccardo Corvi, Matthias Nießner, and Luisa Verdoliva. Raising the bar of ai-generated image detection with clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4356–4366, 2024. 7
- [19] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*, 2023. 2
- [20] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023. 3
- [21] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 1
- [22] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 1, 7, 2
- [23] Guian Fang, Wenbiao Yan, Yuanfan Guo, Jianhua Han, Zutao Jiang, Hang Xu, Shengcai Liao, and Xiaodan Liang.

- Humanrefiner: Benchmarking abnormal human generation and refining with coarse-to-fine pose-reversible guidance. In *European Conference on Computer Vision*, pages 201–217. Springer, 2024. 4, 1
- [24] Joel Frank, Thorsten Eisenhofer, Lea Schönherr, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. Leveraging frequency analysis for deep fake image recognition. In *International conference on machine learning*, pages 3247–3258. PMLR, 2020. 3, 2
- [25] Jessica Fridrich and Jan Kodovsky. Rich models for steganalysis of digital images. *IEEE Transactions on information Forensics and Security*, 7(3):868–882, 2012. 3
- [26] Xingyu Fu, Muyu He, Yujie Lu, William Yang Wang, and Dan Roth. Commonsense-t2i challenge: Can text-to-image generation models understand commonsense? *arXiv preprint arXiv:2406.07546*, 2024. 4, 1, 2
- [27] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27, 2014. 1
- [28] Xiao Guo, Xiufeng Song, Yue Zhang, Xiaohong Liu, and Xiaoming Liu. Rethinking vision-language model in face forensics: Multi-modal interpretable forged face detector. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 105–116, 2025. 1
- [29] Jian Han, Jinlai Liu, Yi Jiang, Bin Yan, Yuqi Zhang, Zehuan Yuan, Bingyue Peng, and Xiaobing Liu. Infinity: Scaling bitwise autoregressive modeling for high-resolution image synthesis, 2024. 7, 2
- [30] Yan Hong and Jianfu Zhang. Wildfake: A large-scale challenging dataset for ai-generated images detection. *arXiv preprint arXiv:2402.11843*, 2024. 3
- [31] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 6
- [32] Zhenglin Huang, Jinwei Hu, Xiangtai Li, Yiwei He, Xingyu Zhao, Bei Peng, Baoyuan Wu, Xiaowei Huang, and Guangliang Cheng. Sida: Social media image deepfake detection, localization and explanation with large multimodal model. *arXiv preprint arXiv:2412.04292*, 2024. 1
- [33] Zhengchao Huang, Bin Xia, Zicheng Lin, Zhun Mou, and Wenming Yang. Ffaa: Multimodal large language model based explainable open-world face forgery analysis assistant. *arXiv preprint arXiv:2408.10072*, 2024. 1, 3, 4, 5
- [34] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 3
- [35] Shan Jia, Reilin Lyu, Kangran Zhao, Yize Chen, Zhiyuan Yan, Yan Ju, Chuanbo Hu, Xin Li, Baoyuan Wu, and Siwei Lyu. Can chatgpt detect deepfakes? a study of using multimodal large language models for media forensics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4324–4333, 2024. 3
- [36] Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L’elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b. *ArXiv*, abs/2310.06825, 2023. 3
- [37] Yan Ju, Shan Jia, Lipeng Ke, Hongfei Xue, Koki Nagano, and Siwei Lyu. Fusing global and local features for generalized ai-synthesized image detection. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 3465–3469. IEEE, 2022. 2
- [38] Negar Kamali, Karyn Nakamura, Angelos Chatzimpampas, Jessica Hullman, and Matthew Groh. How to distinguish ai-generated images from authentic photographs. *arXiv preprint arXiv:2406.08651*, 2024. 4, 5
- [39] Negar Kamali, Karyn Nakamura, Aakriti Kumar, Angelos Chatzimpampas, Jessica Hullman, and Matthew Groh. Characterizing photorealism and artifacts in diffusion model-generated images. *arXiv preprint arXiv:2502.11989*, 2025. 4, 5
- [40] Hengrui Kang, Siwei Wen, Zichen Wen, Junyan Ye, Weijia Li, Peilin Feng, Baichuan Zhou, Bin Wang, Dahua Lin, Linfeng Zhang, et al. Legion: Learning to ground and explain for synthetic image detection. *arXiv preprint arXiv:2503.15264*, 2025. 1
- [41] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 3
- [42] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 1
- [43] Mamadou Keita, Wassim Hamidouche, Hessen Bougueffa Eutamene, Abdenour Hadid, and Abdelmalik Taleb-Ahmed. Bi-lora: A vision-language approach for synthetic image detection. *arXiv preprint arXiv:2404.01959*, 2024. 1, 3
- [44] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 3
- [45] Christos Koutlis and Symeon Papadopoulos. Leveraging representations from intermediate encoder-blocks for synthetic image detection. In *European Conference on Computer Vision*, pages 394–411. Springer, 2024. 1, 7, 2, 4
- [46] Rohit Kundu, Athula Balachandran, and Amit K Roy-Chowdhury. Truthlens: Explainable deepfake detection for face manipulated and fully synthetic data. *arXiv preprint arXiv:2503.15867*, 2025. 1
- [47] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023. 2
- [48] Binh M Le and Simon S Woo. Quality-agnostic deepfake detection with intra-model collaborative learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22378–22389, 2023. 2

- [49] Jiawei Li, Fanrui Zhang, Jiaying Zhu, Esther Sun, Qiang Zhang, and Zheng-Jun Zha. Forgerygpt: Multimodal large language model for explainable image forgery detection and localization. *arXiv preprint arXiv:2410.10238*, 2024. 1
- [50] Ouxiang Li, Jiayin Cai, Yanbin Hao, Xiaolong Jiang, Yao Hu, and Fuli Feng. Improving synthetic image detection towards generalization: An image transformation perspective. *arXiv preprint arXiv:2408.06741*, 2024. 2
- [51] Yixuan Li, Xuelin Liu, Xiaoyang Wang, Shiqi Wang, and Weisi Lin. Fakebench: Uncover the achilles’ heels of fake images with large multimodal models. *arXiv preprint arXiv:2404.13306*, 2024. 2, 3, 4
- [52] Jingchun Lian, Lingyu Liu, Yaxiong Wang, Yujiao Wu, Li Zhu, and Zhedong Zheng. A large-scale interpretable multi-modality benchmark for facial image forgery localization. *arXiv preprint arXiv:2412.19685*, 2024. 1
- [53] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004. 5
- [54] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 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, 2014. 7
- [55] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024. 5, 1, 2, 10, 15
- [56] Bo Liu, Fan Yang, Xiuli Bi, Bin Xiao, Weisheng Li, and Xinbo Gao. Detecting generated images by real images. In *European Conference on Computer Vision*, pages 95–110. Springer, 2022. 2
- [57] Huan Liu, Zichang Tan, Chuangchuang Tan, Yunchao Wei, Yao Zhao, and Jingdong Wang. Forgery-aware adaptive transformer for generalizable synthetic image detection. *arXiv preprint arXiv:2312.16649*, 2023. 3
- [58] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 7
- [59] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 3, 5, 1
- [60] Huan Liu, Zichang Tan, Chuangchuang Tan, Yunchao Wei, Jingdong Wang, and Yao Zhao. Forgery-aware adaptive transformer for generalizable synthetic image detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2
- [61] Yunpeng Luo, Junlong Du, Ke Yan, and Shouhong Ding. Lare²: Latent reconstruction error based method for diffusion-generated image detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17006–17015, 2024. 1, 7, 2
- [62] Fanqing Meng, Wenqi Shao, Lixin Luo, Yahong Wang, Yiran Chen, Quanfeng Lu, Yue Yang, Tianshuo Yang, Kaipeng Zhang, Yu Qiao, et al. Phybench: A physical common-sense benchmark for evaluating text-to-image models. *arXiv preprint arXiv:2406.11802*, 2024. 4, 1, 2
- [63] Mistral-AI. Pixtral large. <https://mistral.ai/news/pixtral-large/>, 2024. 4
- [64] Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. Towards universal fake image detectors that generalize across generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24480–24489, 2023. 3, 2
- [65] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 5
- [66] Siran Peng, Zipei Wang, Li Gao, Xiangyu Zhu, Tianshuo Zhang, Ajian Liu, Haoyuan Zhang, and Zhen Lei. Mllm-enhanced face forgery detection: A vision-language fusion solution. *arXiv preprint arXiv:2505.02013*, 2025. 1
- [67] Chenfan Qu, Jian Liu, Haoxing Chen, Baihan Yu, Jingjing Liu, Weiqiang Wang, and Lianwen Jin. Explainable tampered text detection via multimodal large models. *arXiv preprint arXiv:2412.14816*, 2024. 1, 3, 4
- [68] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3
- [69] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 2024. 2, 3
- [70] Jonas Ricker, Simon Damm, Thorsten Holz, and Asja Fischer. Towards the detection of diffusion model deepfakes. *arXiv preprint arXiv:2210.14571*, 2022. 2
- [71] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 1
- [72] Ayush Sarkar, Hanlin Mai, Amitabh Mahapatra, Svetlana Lazebnik, D.A. Forsyth, and Anand Bhattad. Shadows don’t lie and lines can’t bend! generative models don’t know projective geometry...for now. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 28140–28149, 2024. 4, 1
- [73] Zeyang Sha, Zheng Li, Ning Yu, and Yang Zhang. De-fake: Detection and attribution of fake images generated by text-to-image diffusion models. *arXiv preprint arXiv:2210.06998*, 2022. 2
- [74] Xiufeng Song, Xiao Guo, Jiache Zhang, Qirui Li, Lei Bai, Xiaoming Liu, Guangtao Zhai, and Xiaohong Liu. On learning multi-modal forgery representation for diffusion generated video detection. In *Proceeding of Thirty-eighth Conference on Neural Information Processing Systems*, Vancouver, Canada, 2024. 1
- [75] Ke Sun, Hong Liu, Qixiang Ye, Yue Gao, Jianzhuang Liu, Ling Shao, and Rongrong Ji. Domain general face forgery

- detection by learning to weight. In *Proceedings of the AAAI conference on artificial intelligence*, pages 2638–2646, 2021. [1](#)
- [76] Ke Sun, Hong Liu, Taiping Yao, Xiaoshuai Sun, Shen Chen, Shouhong Ding, and Rongrong Ji. An information theoretic approach for attention-driven face forgery detection. In *European conference on computer vision*, pages 111–127. Springer, 2022.
- [77] Ke Sun, Taiping Yao, Shen Chen, Shouhong Ding, Jilin Li, and Rongrong Ji. Dual contrastive learning for general face forgery detection. In *Proceedings of the AAAI conference on artificial intelligence*, pages 2316–2324, 2022.
- [78] Ke Sun, Shen Chen, Taiping Yao, Xiaoshuai Sun, Shouhong Ding, and Rongrong Ji. Continual face forgery detection via historical distribution preserving. *International Journal of Computer Vision*, 133(3):1067–1084, 2025. [1](#)
- [79] Ke Sun, Shen Chen, Taiping Yao, Ziyin Zhou, Jiayi Ji, Xiaoshuai Sun, Chia-Wen Lin, and Rongrong Ji. Towards general visual-linguistic face forgery detection. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19576–19586, 2025. [1](#)
- [80] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024. [7](#), [2](#)
- [81] Zhihao Sun, Haoran Jiang, Haoran Chen, Yixin Cao, Xipeng Qiu, Zuxuan Wu, and Yu-Gang Jiang. Forgerysluth: Empowering multimodal large language models for image manipulation detection. *arXiv preprint arXiv:2411.19466*, 2024. [1](#), [3](#), [4](#), [5](#)
- [82] Chuangchuang Tan, Yao Zhao, Shikui Wei, Guanhua Gu, and Yunchao Wei. Learning on gradients: Generalized artifacts representation for gan-generated images detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12105–12114, 2023. [2](#)
- [83] Chuangchuang Tan, Yao Zhao, Shikui Wei, Guanhua Gu, Ping Liu, and Yunchao Wei. Rethinking the up-sampling operations in cnn-based generative network for generalizable deepfake detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28130–28139, 2024. [1](#), [3](#), [5](#), [6](#), [7](#), [2](#)
- [84] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. 2024. [7](#), [2](#)
- [85] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *arXiv preprint arXiv:2404.02905*, 2024. [1](#)
- [86] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015. [5](#)
- [87] 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. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. [4](#)
- [88] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot...for now. In *CVPR*, 2020. [2](#), [3](#)
- [89] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8695–8704, 2020. [1](#), [3](#), [7](#), [2](#)
- [90] Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023. [2](#)
- [91] Zhendong Wang, Jianmin Bao, Wengang Zhou, Weilun Wang, Hezhen Hu, Hong Chen, and Houqiang Li. Dire for diffusion-generated image detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22445–22455, 2023. [1](#), [3](#), [2](#)
- [92] Siwei Wen, Junyan Ye, Peilin Feng, Hengrui Kang, Zichen Wen, Yize Chen, Jiang Wu, Wenjun Wu, Conghui He, and Weijia Li. Spot the fake: Large multimodal model-based synthetic image detection with artifact explanation. *arXiv preprint arXiv:2503.14905*, 2025. [1](#)
- [93] Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, Chong Ruan, et al. Janus: Decoupling visual encoding for unified multimodal understanding and generation. *arXiv preprint arXiv:2410.13848*, 2024. [7](#), [2](#)
- [94] Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li, Wenxiu Sun, Qiong Yan, Guangtao Zhai, et al. Q-bench: A benchmark for general-purpose foundation models on low-level vision. *arXiv preprint arXiv:2309.14181*, 2023. [2](#), [5](#)
- [95] H. Wu, J. Zhou, and S. Zhang. Generalizable synthetic image detection via language-guided contrastive learning. *arXiv preprint:2305.13800*, 2023. [2](#)
- [96] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. [7](#), [2](#)
- [97] Zhipei Xu, Xuanyu Zhang, Runyi Li, Zecheng Tang, Qing Huang, and Jian Zhang. Fakeshield: Explainable image forgery detection and localization via multi-modal large language models. *arXiv preprint arXiv:2410.02761*, 2024. [1](#), [3](#), [4](#), [5](#)
- [98] Shilin Yan, Ouxiang Li, Jiayin Cai, Yanbin Hao, Xiaolong Jiang, Yao Hu, and Weidi Xie. A sanity check for ai-generated image detection. *arXiv preprint arXiv:2406.19435*, 2024. [1](#), [3](#), [7](#), [2](#), [4](#)
- [99] Fan Yang, Ru Zhen, Jianing Wang, Yanhao Zhang, Haoxiang Chen, Haonan Lu, Sicheng Zhao, and Guiguang Ding. Heie: Mllm-based hierarchical explainable aigc image implausibility evaluator. *arXiv preprint arXiv:2411.17261*, 2024. [1](#)
- [100] Junyan Ye, Baichuan Zhou, Zilong Huang, Junan Zhang, Tianyi Bai, Hengrui Kang, Jun He, Honglin Lin, Zihao

- Wang, Tong Wu, et al. Loki: A comprehensive synthetic data detection benchmark using large multimodal models. *arXiv preprint arXiv:2410.09732*, 2024. [2](#), [3](#), [4](#)
- [101] Peipeng Yu, Jianwei Fei, Hui Gao, Xuan Feng, Zhihua Xia, and Chip Hong Chang. Unlocking the capabilities of large vision-language models for generalizable and explainable deepfake detection. *arXiv preprint arXiv:2503.14853*, 2025. [1](#)
- [102] Yue Zhang, Ben Colman, Xiao Guo, Ali Shahriyari, and Gaurav Bharaj. Common sense reasoning for deepfake detection. In *European Conference on Computer Vision*, pages 399–415. Springer, 2024. [1](#)
- [103] Yuhui Zhang, Alyssa Unell, Xiaohan Wang, Dhruba Ghosh, Yuchang Su, Ludwig Schmidt, and Serena Yeung-Levy. Why are visually-grounded language models bad at image classification? *arXiv preprint arXiv:2405.18415*, 2024. [2](#), [5](#)
- [104] Yue Zhang, Ben Colman, Xiao Guo, Ali Shahriyari, and Gaurav Bharaj. Common sense reasoning for deepfake detection. In *European Conference on Computer Vision*, pages 399–415. Springer, 2025. [1](#), [3](#)
- [105] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023. [3](#)
- [106] Nan Zhong, Yiran Xu, Zhenxing Qian, and Xinpeng Zhang. Rich and poor texture contrast: A simple yet effective approach for ai-generated image detection. *arXiv preprint arXiv:2311.12397*, 2023. [7](#), [2](#)
- [107] Nan Zhong, Yiran Xu, Sheng Li, Zhenxing Qian, and Xinpeng Zhang. Patchcraft: Exploring texture patch for efficient ai-generated image detection. *arXiv preprint arXiv:2311.12397*, pages 1–18, 2024. [3](#)
- [108] Mingjian Zhu, Hanting Chen, Qiangyu Yan, Xudong Huang, Guanyu Lin, Wei Li, Zhijun Tu, Hailin Hu, Jie Hu, and Yunhe Wang. Genimage: A million-scale benchmark for detecting ai-generated image, 2023. [2](#), [3](#), [4](#)



AIGI-Holmes: Towards Explainable and Generalizable AI-Generated Image Detection via Multimodal Large Language Models

Supplementary Material

A. Prompts

We present the prompts used in our paper as follows: the prompt for querying real images is shown in Fig. 6, and the prompt for querying AI-generated images is shown in Fig. 7. The prompts for validation and evaluation in the Multi-Expert Jury are displayed in Fig. 8. Prompts for filtering images with specific defects and their corresponding queries are shown in Fig. 9, Fig. 10, Fig. 11, Fig. 12, Fig. 13. The prompt for modifying DeepseekV3 [55] based on feedback is shown in Fig. 14.

B. Additional Related Work

In the Deepfake Detection, in addition to existing black-box binary classification methods [11, 14, 75–78], a new wave of explainable detection approaches based on multimodal large models has emerged [13, 28, 46, 52, 66, 79, 101, 102]. X^2 DFD [13] introduced expert-level features for better explanatory outputs, while ForgeryTalker [52] developed a framework for DeepFake facial image localization and explanation using a multimodal forgery tracking dataset (MMTT) and a forgery prompt network (FPN). In Image Forgery Detection, ForgeryGPT [49] integrated a Mask-Aware Forgery Extractor, improving detection and localization tasks (IFDL) with interpretable reasoning. SIDA [32] proposed a framework for detecting deepfakes in social media images, which was trained on a dataset of 300,000 images (SID-Set), including 3,000 images with annotations. In AI-Generated Image Evaluation, HEIE [99] introduced a CoT-driven evaluator and an irrationality mapper, combining heatmaps, scores, and explanations for fine-grained defect localization, supported by the ExplAIGI-Eval dataset. For AI-generated video detection, MM-Det [74] extracted multimodal features from videos using LLaVA [59] and designed a sophisticated fusion strategy. Notably, two contemporaneous works on explainable AI-generated image detection, LEGION [40] and FakeVLM [92], have attracted our attention. We briefly compare these works with ours as follows: In terms of data construction, FakeVLM utilizes MLLM-generated explanations, while LEGION relies solely on human annotation. Our approach leverages MLLM-generated annotations with expert-guided filtering and human corrections, striking a balance between annotation cost and data quality. Regarding training, unlike FakeVLM and LEGION, which only employ supervised fine-tuning (SFT), we introduce a systematic three-stage training pipeline—visual-expert pre-training, SFT, and direct preference optimization (DPO)—to provide

a more comprehensive exploration of adapting MLLMs for explainable AI-generated image detection.

C. Holmes-Set Analysis

In the main text and Fig. 2, we introduced the specific construction process of the Holmes-Set. Furthermore, we present the pseudocode of the Holmes-Set construction process in Algorithm 1 for reference. In the following sections, we will analyze the dataset and visualize some examples.

Expert-guided Filter Method. In the main text, we describe the process of filtering generated images that contain common AI-generated flaws. Below are the expert models we used:

- **Face:** We used a face detection method from OpenCV2 to filter out images containing faces from AI-generated images. MLLMs were also used for annotation and verification, as shown in Fig. 9.
- **Human Body:** We obtained images with various body anomalies from AbHuman [23] and generated red bounding boxes based on the provided bounding boxes to serve as visual prompts for MLLMs, as depicted in Fig. 10.
- **Text & Logos:** We utilized a model from PaddleOCR¹ to filter out images containing text. Additionally, we employed MLLMs for annotation and verification, as illustrated in Fig. 11.
- **Projective Geometry:** We used the model from [72] to filter out images that might contain projective geometry errors. MLLMs were employed for annotation and verification, as illustrated in Fig. 12.
- **Common Sense & Physical Laws:** We generated a batch of images using prompts that induce hallucinations in T2I generation models, as mentioned in commonsense [26] and physical laws [62]. MLLMs were used for annotation and verification, as shown in Fig. 13.

For each type of flaw mentioned above, we ultimately obtained 2,000 annotated image-text pairs, and selected 2,000 images from the COCO dataset as the corresponding real image dataset.

Multi-Expert Jury. We conducted an analysis of the cross-validation scores of the models. In Fig. 16, the leftmost image shows the distribution of scores given by different models to other models, where ModelA_ModelB represents the score given by Model A to Model B. The middle image represents the correlation of scores given to different models, and the rightmost image represents the adoption of annotations from various models in the image-text pairs. From the

Algorithm 1 Holmes-Set Construction

Require: Source datasets \mathcal{M} , Expert datasets \mathcal{E} , MLLM jury \mathcal{J} , SFT budget K , DPO rounds R_1, R_2

```
1: // Stage 1: Data Collection
2:  $\mathcal{D}_{\text{base}} \leftarrow \text{Select}(\mathcal{M}, 45\text{K})$   $\triangleright$  From CNNDetection, GenImage, etc.
3:  $\mathcal{D}_{\text{expert}} \leftarrow \text{Filter}(\mathcal{E}, 20\text{K})$   $\triangleright$  Expert-guided flaws
4:  $\mathcal{D}_{\text{gen}} \leftarrow \text{Generate}([26, 62])$   $\triangleright$  Common sense/physics flaws
5:  $\mathcal{D}_{\text{all}} \leftarrow \mathcal{D}_{\text{base}} \cup \mathcal{D}_{\text{expert}} \cup \mathcal{D}_{\text{gen}}$ 

6: // Stage 2: Automated Annotation
7:  $\mathcal{D}_{\text{SFT}}, \mathcal{D}_1 \leftarrow \emptyset$ 
8: for each image  $I \in \mathcal{D}_{\text{all}}$  do
9:   if  $I \in \mathcal{D}_{\text{expert}} \cup \mathcal{D}_{\text{gen}}$  then
10:     $x_p \leftarrow \text{SpecialistPrompt}(I)$   $\triangleright$  Focus on known flaws
11:   else
12:     $x_p \leftarrow \text{GeneralPositivePrompt}()$   $\triangleright$  Judge criteria
13:     $x_n \leftarrow \text{GeneralPositivePrompt}()$ 
14:   end if
15:    $A \leftarrow \mathcal{J}(I, x_p)$   $\triangleright$  Jury system generates annotation
16:    $A_- \leftarrow \mathcal{J}(I, x_n)$   $\triangleright$  Jury system generates annotation

17: // Self-verification for SFTSet
18:  $\text{score} \leftarrow \frac{1}{|\mathcal{J}|} \sum_{J \in \mathcal{J}} \text{Score}(A, J(I))$ 
19: if  $\text{score} \geq \delta_{\text{SFT}}$  then
20:    $\mathcal{D}_{\text{SFT}} \leftarrow \mathcal{D}_{\text{SFT}} \cup \{(I, A)\}$ 
21: end if
22:  $\mathcal{D}_1 \leftarrow \mathcal{D}_1 \cup \{(I, A, A_-)\}$ 
23: end for

24: // Stage 3: Human Preference Refinement
25:  $\mathcal{D}_2 \leftarrow \emptyset$ 
26: Sample  $\mathcal{D}_{\text{human}} \subset \mathcal{M} \cup \mathcal{E}$  (2K images)
27: Sample  $\mathcal{D}_{\text{mllm}} \subset \mathcal{M} \cup \mathcal{E}$  (2K images)
28: for each  $(I, A) \in \mathcal{D}_{\text{human}}$  do
29:    $A'_1 \leftarrow \text{HumanRevise}(A)$   $\triangleright$  Human-Expert modifications
30:    $A' \leftarrow \text{DeepseekV3}(A, A'_1)$   $\triangleright$  LLM-assisted refinement
31:    $\mathcal{D}_2 \leftarrow \mathcal{D}_2 \cup \{(I, A')\}$ 
32: end for
33: for each  $(I, A) \in \mathcal{D}_{\text{mllm}}$  do
34:    $A'_2 \leftarrow \text{MLLMRevise}(A)$   $\triangleright$  MLLM-Expert modifications
35:    $A'' \leftarrow \text{DeepseekV3}(A, A'_2)$   $\triangleright$  LLM-assisted refinement
36:    $\mathcal{D}_2 \leftarrow \mathcal{D}_2 \cup \{(I, A'')\}$ 
37: end for

38: // Stage 4: Comprehensive Evaluation
39:  $\mathcal{D}_{\text{test}} \leftarrow \text{HumanRevise}(1\text{K})$ 
40: Compute  $\{\text{BLEU}, \text{CIDEr}, \text{ROUGE-L}, \text{METER}\}$  on  $\mathcal{D}_{\text{test}}$ 
41:  $\text{MLLM-Score} \leftarrow \frac{1}{|\mathcal{J}|} \sum_{J \in \mathcal{J}} \text{BatchScore}(\mathcal{D}_{\text{test}}, J)$ 
42:  $\text{Human-Score} \leftarrow \text{ExpertPreferenceAssessment}(\mathcal{D}_{\text{test}})$ 

43: return  $\mathcal{D}_{\text{SFT}}, \mathcal{D}_1, \mathcal{D}_2$ 
```

figures, we can observe that the scoring differentiation of the Pixtral-124B and InternVL-76B models is relatively higher,

while the scoring correlation between InternVL2.5-78B and Qwen2VL-72B is very high. This may be due to the use of the same large language model base, Qwen, which could have affected these models' ability to evaluate AI-generated image explanations due to safety alignment operations. Ultimately, our dataset is composed mostly of annotations from Pixtral-124B, with some from InternVL-76B, Qwen2VL-72B, and a small amount from InternVL2.5-78B. We present some examples from our SFT dataset in Fig. 21, Fig. 22, Fig. 23, and Fig. 24. For InternVL-76B, Qwen2VL-72B, and InternVL2.5-78B, we utilized the vllm framework [47], which is widely used for large model inference, to perform local deployment. For Pixtral-124B, we accessed it via the official website's API.

Holmes-DPOSet. In the main text, we introduced how we solicited revision suggestions for the responses of the supervised fine-tuning model from human experts and multimodal large language model experts. These suggestions were then fed into DeepseekV3 [55] using the prompt shown in Fig. 14 to obtain the final human-aligned preference samples. This process is illustrated in Fig. 20. The revision suggestions originated from the annotation platform we developed, as shown in Fig. 19, based on LabelLLM². By referring to the revision suggestions, the model effectively supplemented the original response by addressing an additional anatomical error that was previously overlooked, while maintaining a high level of consistency before and after the modification. This is beneficial for executing the DPO training phase. In Fig. 25 and Fig. 26, we present sample pairs from \mathcal{D}_1 and \mathcal{D}_2 , respectively.

Construction Process of \mathcal{P}_3 . For models such as Janus [93], Janus-Pro [12], VAR [84], Infinity [29], Show-o [96], LlamaGen [80], FLUX [4], and SD3.5 [22], we deployed these models to perform inference and generate images. The test images for PixArt-XL [9] and some of the FLUX [4] images were sourced from [50]. The image resolutions are as follows: FLUX, SD3.5, Infinity, and PixArt-XL have a resolution of 1024×1024 ; Show-o and LlamaGen have a resolution of 512×512 ; Janus and Janus-Pro have a resolution of 384×384 ; and VAR has a resolution of 256×256 . Examples of these images can be found in Fig. 15.

D. More Experiments

D.1. ALL Baselines

We compared various baselines including CNNSpot [89], FreDect [24], Fusing [37], LNP [56], LGrad [82], UnivFD [64], DIRE [91], PatchCraft [106], NPR [83], AntiFakePrompt [5], Fatformer [60], Ricker2022 [70], DEFAKE [73], LASTED [95], QAD [48], InstructBLIP [19], CogVLM [90], LaRE [61], RINE [45], AIDE [98]. These

¹<https://github.com/PaddlePaddle/PaddleOCR>

²<https://github.com/opendatalab/LabelLLM>

baseline methods are trained and tested under fair and consistent conditions across different settings.

D.2. Training Details of Baselines in Protocol-III

The training details for the methods we trained on the Holmes-SFTSet training set are as follows, with all unmentioned details consistent with the code provided by the original authors:

- **CNNSpot**: We trained using the Adam optimizer with a learning rate of $2e-4$ and a batch size of 32 for 10 epochs.
- **AntifakePrompt**: We trained using the Adam optimizer with a learning rate of $2e-4$ and a batch size of 8 for 10 epochs.
- **UnivFD**: We trained using the Adam optimizer with a learning rate of $1e-4$ and a batch size of 48 for 20 epochs.
- **NPR**: We trained using the Adam optimizer with a learning rate of $2e-4$ and a batch size of 32 for 30 epochs.
- **LaRE**: We trained using the Adam optimizer with a learning rate of $1e-4$ and a batch size of 48 for 10 epochs.
- **RINE**: We trained using the Adam optimizer with a learning rate of $2e-4$ and a batch size of 16 for 5 epochs.
- **AIDE**: We trained using the Adam optimizer with a learning rate of $2e-4$ and a batch size of 8 for 5 epochs.

D.3. Comparisons with SOTAs

Protocol-I. The quantitative results in Tab. 6 show the classification accuracy of various methods and generators within the range of \mathcal{P}_1 . In this evaluation, all methods were trained solely on the four categories (car, cat, chair, horse) generated by ProGAN, except for DIRE-D, which was trained on the Diffusion dataset of ADM. AIGI-Holmes demonstrates significant improvements compared to the current state-of-the-art (SOTA) methods PatchCraft and AIDE, with average accuracy increases of 3.9% and 0.4%, respectively. The AIDE method integrates semantic, low-frequency, and high-frequency information through a dual-stream structure, showing remarkable detection performance for some Diffusion methods. However, our method not only performs well in detecting Diffusion methods but also excels in GauGAN and BigGAN methods, where these SOTA methods underperform, with improvements of 10.6% and 21.8%, respectively. Although trained only on the single forgery method of ProGAN, which can easily lead to overfitting for mllm, our model still demonstrates good generalization to unseen diffusion methods, highlighting the potential of our approach.

Protocol-II. The quantitative results in Tab. 7 show the classification accuracy of various methods and generators within the range of \mathcal{P}_2 . All methods were trained or fine-tuned on the dataset primarily generated by SD3, highlighted in gray in the table. AIGI-Holmes shows significant improvements compared to the current state-of-the-art (SOTA) method AntifakePrompt, with an average accuracy increase of 1.14%. The AntifakePrompt method uses prompt learning to per-

form binary classification on images with mllm, achieving excellent detection performance on a large number of unseen diffusion methods, image inpainting methods, image super-resolution methods, and Deepfake datasets. However, our method achieves over 10% improvement in detection accuracy on LaMa, DALLÉ-3, and real image datasets COCO and Flickr. In addition to providing accurate binary classification results, it can also output explanations and reasons corresponding to the predicted results.

D.4. More Ablation Study

We primarily conduct ablation experiments on \mathcal{P}_1 and \mathcal{P}_3 to evaluate the prediction accuracy of the models.

Large language models. We utilize the mainstream multimodal architecture LLaVA [59] and conduct ablation experiments on this architecture using large language models. Specifically, we experiment with three language models as shown in Tab. 8: Llama3-8B [2], Mistral-7B [36], and Vicuna-7B [105]. The results indicate that the impact on accuracy is minimal when using different large language models, with Mistral-7B showing a slight advantage. Consequently, our final approach employs Mistral-7B as the large language model within the LLaVA architecture.

Training methods. We conduct ablation experiments on different training approaches to demonstrate the necessity and effectiveness of our Holmes Pipeline. Specifically, we investigate two training strategies: fine-tuning only the large language model using LoRA, and simultaneously applying LoRA training to both the visual component and the large language model. The results, presented in Tab. 9, indicate that our Holmes Pipeline effectively adapts the multimodal architecture to the task of AI-generated image detection. Compared to alternative approaches, our method achieves significant improvements in accuracy, with enhancements of 2.0% and 7.0% on \mathcal{P}_1 and \mathcal{P}_3 , respectively.

Integration of Visual Experts. We conduct ablation studies on the visual expert components introduced during collaborative decoding. The experimental results, presented in Tab. 10, demonstrate that both types of visual experts enhance the performance of the method on \mathcal{P}_1 and \mathcal{P}_3 . The improvement on \mathcal{P}_1 is more pronounced, with the integration of a single visual expert and all visual experts increasing the detection accuracy by 5.2% and 6.4%, respectively. For \mathcal{P}_3 , the improvements are 0.7% and 1.0%, respectively. This discrepancy may be attributed to the fact that \mathcal{P}_1 contains only one type of forgery, namely Progan, leading to overfitting in the MLLM and limiting its generalizability. In contrast, the visual experts exhibit better generalization capabilities. Through collaborative decoding, we ensure the generalizability of the model’s detection performance.

Impact of Dataset Size. Due to the high cost of data annotation, we used a relatively small training set ($\sim 65K$) to train the baselines, which may not be sufficient for the baselines

Method	<i>ProGAN</i>	<i>StyleGAN</i>	<i>BigGAN</i>	<i>CycleGAN</i>	<i>StarGAN</i>	<i>GauGAN</i>	<i>StyleGAN2</i>	<i>WFR</i>	<i>ADM</i>	<i>Glide</i>	<i>Midjourney</i>	<i>SD v1.4</i>	<i>SD v1.5</i>	<i>VQDM</i>	<i>Wukong</i>	<i>DALLE2</i>	<i>Mean</i>
CNNSpot	100.00	90.17	71.17	87.62	94.60	81.42	86.91	91.65	60.39	58.07	51.39	50.57	50.53	56.46	51.03	50.45	70.78
FreDect	99.36	78.02	81.97	78.77	94.62	80.57	66.19	50.75	63.42	54.13	45.87	38.79	39.21	77.80	40.30	34.70	64.03
Fusing	100.00	85.20	77.40	87.00	97.00	77.00	83.30	66.80	49.00	57.20	52.20	51.00	51.40	55.10	51.70	52.80	68.38
LNP	99.67	91.75	77.75	84.10	99.92	75.39	94.64	70.85	84.73	80.52	65.55	85.55	85.67	74.46	82.06	88.75	83.84
LGrad	99.83	91.08	85.62	86.94	99.27	78.46	85.32	55.70	67.15	66.11	65.35	63.02	63.67	72.99	59.55	65.45	75.34
UnivFD	99.81	84.93	<u>95.08</u>	<u>98.33</u>	95.75	99.47	74.96	86.90	66.87	62.46	56.13	63.66	63.49	85.31	70.93	50.75	78.43
DIRE-G	95.19	83.03	70.12	74.19	95.47	67.79	75.31	58.05	75.78	71.75	58.01	49.74	49.83	53.68	54.46	66.48	68.68
DIRE-D	52.75	51.31	49.70	49.58	46.72	51.23	51.72	53.30	98.25	<u>92.42</u>	<u>89.45</u>	91.24	91.63	<u>91.90</u>	90.90	<u>92.45</u>	71.53
PatchCraft	100.00	92.77	95.80	70.17	99.97	71.58	89.55	85.80	82.17	83.79	90.12	95.38	95.30	88.91	<u>91.07</u>	96.60	<u>89.31</u>
NPR	99.79	97.70	84.35	96.10	99.35	82.50	<u>98.38</u>	65.80	69.69	78.36	77.85	78.63	78.89	78.13	76.11	64.90	82.91
AIDE	99.99	99.64	83.95	98.48	99.91	73.25	98.00	94.20	<u>93.43</u>	95.09	77.20	<u>93.00</u>	<u>92.85</u>	95.16	93.55	96.60	<u>92.77</u>
AIGI-Holmes	100.00	<u>98.35</u>	94.51	97.03	100.00	<u>95.19</u>	98.88	95.71	88.43	91.53	81.56	91.28	91.38	90.94	89.46	85.32	93.16

Table 6. Evaluation on \mathcal{P}_1 : All baseline results are trained on \mathcal{P}_1 's training set to ensure a fair comparison. The remaining baseline results are sourced from AIDE [98].

Dataset	Ricker2022	ResNet	FatFormer	CNNSpot		DE-FAKE		DIRE		LASTED		QAD		CogVLM	InstructBLIP		AntifakePrompt		Ours
	P	F	P	P	F	P	F	P	F	P	F	P	F	P	P	LoRA	Orig.	+LaMa	+LaMa
COCO	95.60	<u>99.43</u>	97.40	96.87	<u>99.97</u>	85.97	<u>83.30</u>	81.77	<u>99.93</u>	75.47	<u>58.10</u>	59.57	96.83	98.43	98.93	97.63	92.53	90.40	100.00
Flickr	95.80	99.23	98.13	96.67	100.00	90.67	84.38	77.53	<u>99.93</u>	76.33	65.58	60.23	98.30	99.63	99.63	97.50	91.57	90.60	100.00
SD2	81.10	2.50	16.83	0.17	5.23	<u>97.10</u>	<u>88.07</u>	3.83	30.47	58.69	52.53	51.00	10.67	52.47	40.27	89.57	98.33	97.97	89.61
SD3	88.40	<u>99.83</u>	21.50	4.70	<u>8.60</u>	96.50	<u>95.17</u>	0.00	<u>98.53</u>	78.68	<u>79.51</u>	46.53	99.97	2.10	1.47	<u>97.60</u>	<u>96.17</u>	<u>96.10</u>	<u>99.27</u>
SDXL	81.10	0.50	30.39	0.17	1.53	90.50	72.17	18.17	19.73	51.33	77.65	41.60	9.87	32.57	23.07	96.47	<u>99.17</u>	99.37	99.98
IF	92.65	4.40	27.73	19.17	4.93	99.20	95.20	6.93	63.17	57.99	55.63	59.07	15.17	29.03	20.63	87.90	<u>97.10</u>	95.97	96.37
DALLE-2	52.10	12.80	76.03	3.40	0.87	68.97	61.17	2.13	1.50	57.96	81.91	41.70	14.63	60.70	41.77	<u>99.27</u>	97.27	98.00	100.00
DALLE-3	<u>95.20</u>	2.10	43.97	18.17	3.20	89.00	71.57	0.10	36.27	51.83	53.00	51.23	9.83	6.03	6.63	67.87	80.80	82.97	99.68
playground v2.5	94.40	0.20	29.83	15.73	0.47	96.20	86.77	0.17	17.73	70.95	65.42	38.73	2.47	13.37	6.70	95.43	97.73	<u>98.13</u>	98.53
DiffusionDB	81.20	4.69	60.50	9.03	4.50	80.80	78.10	2.53	16.40	86.48	67.42	52.07	12.07	6.05	53.00	85.40	98.47	<u>98.90</u>	99.25
SGXL	100.00	1.63	97.73	79.30	2.13	56.90	50.20	45.27	9.50	64.39	65.59	46.40	4.20	60.40	69.53	91.20	99.03	<u>99.37</u>	89.09
GLIDE	83.80	49.97	79.80	17.23	5.87	76.50	50.20	4.63	41.77	54.46	68.19	53.63	50.27	59.90	37.97	92.63	98.90	<u>99.70</u>	99.81
Stylization	75.50	0.90	85.03	11.40	4.17	63.97	55.17	9.90	6.30	50.70	67.79	51.93	7.93	42.90	33.97	82.80	94.10	<u>95.77</u>	96.03
DF	14.20	34.20	5.10	0.30	0.03	86.97	77.17	0.27	3.77	86.38	59.36	97.43	22.73	13.80	13.83	67.43	95.03	98.40	<u>98.07</u>
DFDC	46.90	14.20	1.60	0.00	0.00	56.13	48.57	60.13	1.03	70.19	72.42	90.40	28.50	9.00	14.07	85.47	<u>99.83</u>	99.93	89.11
FF++	20.30	37.53	71.30	5.23	0.23	78.90	70.63	25.50	31.93	70.69	56.50	99.47	30.77	35.66	44.20	88.30	95.63	<u>97.97</u>	97.80
LaMa	64.30	1.87	67.03	7.53	0.07	13.03	23.00	13.23	19.47	60.53	97.67	42.03	3.80	5.20	10.90	42.73	39.40	55.80	<u>95.40</u>
SD2IP	59.10	<u>99.76</u>	85.07	1.27	<u>7.23</u>	16.00	<u>75.57</u>	11.37	<u>86.40</u>	56.96	99.87	42.73	<u>96.30</u>	35.50	44.23	91.13	<u>80.80</u>	<u>89.03</u>	<u>94.17</u>
LIIF	58.90	94.43	6.60	8.30	1.07	9.73	53.67	1.10	48.77	56.46	87.34	48.07	95.83	23.47	<u>99.93</u>	84.63	98.50	99.97	60.62
SD2SR	73.90	97.79	84.03	1.40	0.13	29.70	96.67	2.77	27.20	59.59	99.73	47.50	8.63	55.06	69.10	99.90	99.43	<u>99.80</u>	99.76
Average	72.79	37.22	51.42	18.12	10.68	68.45	70.22	16.94	36.68	64.69	70.01	55.64	33.39	37.59	41.83	84.17	91.16	<u>93.08</u>	94.22

Table 7. Evaluation on \mathcal{P}_2 : All baseline results are trained on \mathcal{P}_2 's training set to ensure a fair comparison. The remaining baseline results are sourced from AntifakePrompt [5].

to reach their optimal performance. For a fair comparison, we proportionally increased the size of the training set and retrained the baselines. As shown in Tab. 11, the results indicate that the marginal benefit of increasing dataset size for the baselines diminishes. In contrast, our method achieves optimal test performance even when trained on only a quarter of the dataset, highlighting the efficiency of our training procedure.

Cross-benchmark Evaluation. Given the strong performance of our model on \mathcal{P}_3 , we further evaluate its detection capability across multiple benchmarks to comprehensively demonstrate the effectiveness of our training set and pipeline. Specifically, we compare our model with the two best-performing baselines, RINE [45] and AIDE [98], on LOKI [100], Chameleon [98], \mathcal{P}_1 , and \mathcal{P}_2 . On the two challenging benchmarks, LOKI and Chameleon, our model surpasses the second-best model by 3.4% and 9.2%, respectively, demonstrating its impressive detection capability.

D.5. Comparisons with MLLMs

As shown in the main text, we used the pairwise comparison method from our previous work [17] to compare the responses of different MLLMs. We plotted heatmaps of model comparisons, where the numbers in each heatmap represent the number of times Model A's response was preferred over Model B's response. It can be observed that our model's responses generally achieved better human preference results compared to other models. This demonstrates the effectiveness of our Holmes-DPO. The results in Tab. 3 in the main text were obtained using the algorithm presented in Algorithm 2.

E. More Qualitative Results

In Fig. 27, Fig. 28, Fig. 29, and Fig. 30, we present a comparison of the explanations provided by the baseline MLLMs method discussed in the main text. Additionally, we illustrate the differences between the SFT-tuned model AIGI-Holmes

LLM	\mathcal{P}_1	\mathcal{P}_3
Llama3-8B	91.5	97.8
Mistral-7B	93.2	99.2
Vicuna-7B	93.0	98.3

Table 8. Mean accuracy on \mathcal{P}_1 and \mathcal{P}_3 for Different LLMs

Training Method	\mathcal{P}_1	\mathcal{P}_3
Lora(only LLM)	83.3	90.1
ALL Lora	84.8	91.2
Holmes Pipeline	86.8	98.2

Table 9. Mean accuracy on \mathcal{P}_1 and \mathcal{P}_3 for different training methods.

CLIP	NPR	\mathcal{P}_1	\mathcal{P}_3
		86.8	98.2
✓		92.0	98.9
✓	✓	93.2	99.2

Table 10. Mean accuracy on \mathcal{P}_1 and \mathcal{P}_3 for different visual expert

Scaling	CNNSpot	Antifakeprompt	UnivFD	NPR	LaRE	RINE	AIDE	AIGI-Holmes
1×	72.9	83.9	83.6	84.0	85.0	96.2	97.0	99.2
2×	73.6	84.8	83.7	84.7	85.2	96.8	97.3	-
4×	73.9	85.2	83.9	84.9	85.4	97.0	97.4	-

Table 11. Performance comparison across different scaling factors

	LOKI	Chameleon	\mathcal{P}_1	\mathcal{P}_2
RINE	0.790	0.562	0.960	0.903
AIDE	0.706	0.667	0.949	0.945
AIGI-Holmes	0.824	0.759	0.987	0.954

Table 12. Cross-dataset evaluation results

Algorithm 2 The method for calculating ELO ratings.

```

1: Initialize:
2:  $r \leftarrow \text{defaultdict}(\lambda : \text{INIT\_RATING})$ 
3:  $K \leftarrow 4$ 
4:  $\text{SCALE} \leftarrow 400$ 
5:  $\text{BASE} \leftarrow 10$ 
6:  $\text{INIT\_RATING} \leftarrow 1000$ 
7: for each  $key$  in  $\text{dic}$  do
8:    $\text{model\_a} \leftarrow \text{split}(key, "-")[0]$ 
9:    $\text{model\_b} \leftarrow \text{split}(key, "-")[1]$ 
10:   $\text{winner} \leftarrow \text{dic}[key]$ 
11:   $r_a \leftarrow r[\text{model\_a}]$ 
12:   $r_b \leftarrow r[\text{model\_b}]$ 
13:   $e_a \leftarrow \frac{1}{1 + \text{BASE} \cdot \frac{r_b - r_a}{\text{SCALE}}}$ 
14:   $e_b \leftarrow \frac{1}{1 + \text{BASE} \cdot \frac{r_a - r_b}{\text{SCALE}}}$ 
15:  if  $\text{winner} = \text{"choice\_A"}$  then
16:     $s_a \leftarrow 1$ 
17:  else if  $\text{winner} = \text{"choice\_B"}$  then
18:     $s_a \leftarrow 0$ 
19:  else if  $\text{winner} = \text{"choice\_C"}$  or  $\text{winner} = \text{None}$  then
20:     $s_a \leftarrow 0.5$ 
21:  else
22:    raise  $\text{Exception}(\text{"unexpected vote"} \text{ winner})$ 
23:  end if
24:   $r[\text{model\_a}] \leftarrow r[\text{model\_a}] + K \cdot (s_a - e_a)$ 
25:   $r[\text{model\_b}] \leftarrow r[\text{model\_b}] + K \cdot (1 - s_a - e_b)$ 
26: end for

```

(SFT) and the DPO-tuned model AIGI-Holmes (DPO). The AIGI-Holmes (DPO) model demonstrates a higher quality of responses.

F. Limitations and Future Works

We acknowledge two key limitations. First, as generative models rapidly evolve, the types of forgery-related errors may change, potentially reducing the relevance of our current explanation categories. Our proposed method serves as a foundational approach, and we are committed to extending it to accommodate these emerging error types through adaptable interpretative methods. Second, constrained by the current dataset organization, AIGI-Holmes is limited to generating forensic reports and lacks the image-text dialogue capabilities inherent in multimodal large language models. This limitation of report-only output is also noted in related works [13, 33, 81, 97]. Future work will focus on three aspects to address these limitations: (1) Continuously deploying AIGI-Holmes in real-world scenarios to build larger-scale SFT and DPO datasets, enhancing both robustness against evolving forgeries and explanatory capabilities; (2) Unifying data from Image Forgery Detection and DeepFake Detection within AIGI-Holmes to develop a comprehensive image authenticity detection model; (3) Expanding the dataset to a dialogue format, potentially via specialized tokens to isolate capabilities, thereby equipping the model with multimodal dialogue functionality.

You are an AI visual assistant that can help humans analyze images that may have been generated by AI. You will receive a real image. Your task is to uncover clues that suggest the image may be real.

You need to perform the following two tasks:

1. Please describe the image in detail. At the beginning of the description, state, "**This is a real image.**"
2. Identify clues that differentiate this image from AI-generated images and describe them. You can consider the following angles, but are not limited to them, to find signs that the image is real. For each angle, provide a detailed explanation:

(1) **Line segments:** Carefully observe the overall style, color, and details of the image to determine if there are visual inconsistencies. Pay special attention to the consistency of light, shadows, and colors, and whether there are unnatural areas or traces. Additionally, you need to pay attention to other line segments that may appear in the image, and check for any unnatural distortions or traces. If the image contains multiple objects or people, they should be illuminated by the same light source, and shadows should be consistent. Inconsistent lighting and shadows may indicate that the image is AI-generated. Furthermore, an image will have several vanishing points, and parallel lines in the image should converge at these vanishing points. Parallel lines that do not converge may indicate that the image is AI-generated.

(2) **Edges:** Check for unnatural pixel distribution or edges in the image. Pay particular attention to discontinuous or inconsistent edges, as well as obvious AI-generated traces.

(3) **Texture:** Examine the image for unnatural textures, especially blurry or unclear textures or repetitive texture patterns.

(4) **Distortion:** Check if objects in the image exhibit unnatural distortion, especially whether they conform to the shapes of real objects.

(5) **Overall Hue:** Detect whether the overall tone of the image is overly vibrant, too dull, or contains inconsistencies, as these areas often hint at AI generation.

(6) **Clarity:** Check the image for resolution and compression artifacts. AI-generated images may show unnatural pixel blurriness, jagged edges, or excessive compression traces.

(7) **Perspective:** Observe the perspective and proportional relationships in the image. The perspective and proportions in a real photo should be consistent; if there are abnormal or unnatural perspective relationships, it may indicate signs of being AI-generated. Check whether the depth of field changes reasonably, i.e., whether the blurring between the foreground, background, and subject conforms to actual physical laws.

(8) **Shadows:** Observe whether there are reasonable reflections and shadows in the photo. Real photos typically produce reflections and shadows based on light sources, while composite photos may have unnatural shadows or reflections.

(9) **Text:** If the photo contains text or logos, check whether they are clear and readable and consistent with the surrounding environment. AI-generated images often include unnatural, inconsistent, or unreadable text or logos.

(10) **Physical Laws:** Check whether the content of the image violates physical laws.

(11) **Faces:** Check whether the image content includes faces; AI-generated images often exhibit unnatural facial structures. This includes but is not limited to: 1. Overly perfect facial features lacking the imperfections of real faces. 2. Blurred and distorted edges and facial features. AI-generated faces may appear overly rough, leading to blurriness and distortion at the facial edges. Additionally, you can pay attention to the facial features generated by AI, which often have blurry traces and exaggerated, distorted expressions that differ from real-world faces. 3. Exaggerated and unnatural facial expressions. 5. Inconsistent skin color and texture. 6. Makeup. Overly perfect makeup or makeup inconsistent with age and gender are common flaws in AI-generated images. 7. Body accessories and clothing.

(12) **Body Structure:** Check whether the image content includes body structures; generally, AI-generated images will include some anatomical errors. The body parts you need to consider include head, neck, body, arm, hand, leg, and foot.

(13) **Common Sense:** Check whether there are any violations of common sense in the image.

You can respond in the following format:

This is a real image.

Image Description: xxxxx

Based on the provided real image, here are the reasons why this image is real:

1. xxxxx
2. xxxxx
3. xxxxx

Please provide the most relevant clues that indicate this is a real image, avoiding vague or uninformative responses

Figure 6. General Positive Prompt for annotating real images and General Negative Prompt for annotating AI-generated images.

You are an AI visual assistant that can help humans analyze images that may have been generated by AI. You will receive an AI-generated image. Your task is to uncover clues that suggest the image may be AI-generated. You need to perform the following two tasks:

1. Please describe the image in detail. At the beginning of the description, state, "**This is a fake image.**"
2. Identify clues that differentiate this image from real images and describe them.

You can consider the following angles, but are not limited to them, to find signs that the image is AI-generated. For each angle, provide a detailed explanation:

(1) **Line segments:** Carefully observe the overall style, color, and details of the image to determine if there are visual inconsistencies. Pay special attention to the consistency of light, shadows, and colors, and whether there are unnatural areas or traces. Additionally, you need to pay attention to other line segments that may appear in the image, and check for any unnatural distortions or traces. If the image contains multiple objects or people, they should be illuminated by the same light source, and shadows should be consistent. Inconsistent lighting and shadows may indicate that the image is AI-generated. Furthermore, an image will have several vanishing points, and parallel lines in the image should converge at these vanishing points. Parallel lines that do not converge may indicate that the image is AI-generated.

(2) **Edges:** Check for unnatural pixel distribution or edges in the image. Pay particular attention to discontinuous or inconsistent edges, as well as obvious AI-generated traces.

(3) **Texture:** Examine the image for unnatural textures, especially blurry or unclear textures or repetitive texture patterns.

(4) **Distortion:** Check if objects in the image exhibit unnatural distortion, especially whether they conform to the shapes of real objects.

(5) **Overall Hue:** Detect whether the overall tone of the image is overly vibrant, too dull, or contains inconsistencies, as these areas often hint at AI generation.

(6) **Clarity:** Check the image for resolution and compression artifacts. AI-generated images may show unnatural pixel blurriness, jagged edges, or excessive compression traces.

(7) **Perspective:** Observe the perspective and proportional relationships in the image. The perspective and proportions in a real photo should be consistent; if there are abnormal or unnatural perspective relationships, it may indicate signs of being AI-generated. Check whether the depth of field changes reasonably, i.e., whether the blurring between the foreground, background, and subject conforms to actual physical laws.

(8) **Shadows:** Observe whether there are reasonable reflections and shadows in the photo. Real photos typically produce reflections and shadows based on light sources, while composite photos may have unnatural shadows or reflections.

(9) **Text:** If the photo contains text or logos, check whether they are clear and readable and consistent with the surrounding environment. AI-generated images often include unnatural, inconsistent, or unreadable text or logos.

(10) **Physical Laws:** Check whether the content of the image violates physical laws.

(11) **Faces:** Check whether the image content includes faces; AI-generated images often exhibit unnatural facial structures. This includes but is not limited to: 1. Overly perfect facial features lacking the imperfections of real faces. 2. Blurred and distorted edges and facial features. AI-generated faces may appear overly rough, leading to blurriness and distortion at the facial edges. Additionally, you can pay attention to the facial features generated by AI, which often have blurry traces and exaggerated, distorted expressions that differ from real-world faces. 3. Exaggerated and unnatural facial expressions. 5. Inconsistent skin color and texture. 6. Makeup. Overly perfect makeup or makeup inconsistent with age and gender are common flaws in AI-generated images. 7. Body accessories and clothing.

(12) **Body Structure:** Check whether the image content includes body structures; generally, AI-generated images will include some anatomical errors. The body parts you need to consider include head, neck, body, arm, hand, leg, and foot.

(13) **Common Sense:** Check whether there are any violations of common sense in the image.

You can respond in the following format:

This is a fake image.

Image Description: xxxxx

Based on the provided fake image, here are the reasons why this image is fake:

1. xxxxx
2. xxxxx
3. xxxxx

Please provide the most relevant clues that indicate this is a fake image, avoiding vague or uninformative responses.

Figure 7. General Positive Prompt for annotating AI-generated images and General Negative Prompt for annotating real images.

You can consider the following angles, but are not limited to them, to find signs that the image is real. For each angle, provide a detailed explanation:

(1) **Line segments:** Carefully observe the overall style, color, and details of the image to determine if there are visual inconsistencies. Pay special attention to the consistency of light, shadows, and colors, and whether there are unnatural areas or traces. Additionally, you need to pay attention to other line segments that may appear in the image, and check for any unnatural distortions or traces. If the image contains multiple objects or people, they should be illuminated by the same light source, and shadows should be consistent. Inconsistent lighting and shadows may indicate that the image is AI-generated. Furthermore, an image will have several vanishing points, and parallel lines in the image should converge at these vanishing points. Parallel lines that do not converge may indicate that the image is AI-generated.

(2) **Edges:** Check for unnatural pixel distribution or edges in the image. Pay particular attention to discontinuous or inconsistent edges, as well as obvious AI-generated traces.

(3) **Texture:** Examine the image for unnatural textures, especially blurry or unclear textures or repetitive texture patterns.

(4) **Distortion:** Check if objects in the image exhibit unnatural distortion, especially whether they conform to the shapes of real objects.

(5) **Overall Hue:** Detect whether the overall tone of the image is overly vibrant, too dull, or contains inconsistencies, as these areas often hint at AI generation.

(6) **Clarity:** Check the image for resolution and compression artifacts. AI-generated images may show unnatural pixel blurriness, jagged edges, or excessive compression traces.

(7) **Perspective:** Observe the perspective and proportional relationships in the image. The perspective and proportions in a real photo should be consistent; if there are abnormal or unnatural perspective relationships, it may indicate signs of being AI-generated. Check whether the depth of field changes reasonably, i.e., whether the blurring between the foreground, background, and subject conforms to actual physical laws.

(8) **Shadows:** Observe whether there are reasonable reflections and shadows in the photo. Real photos typically produce reflections and shadows based on light sources, while composite photos may have unnatural shadows or reflections.

(9) **Text:** If the photo contains text or logos, check whether they are clear and readable and consistent with the surrounding environment. AI-generated images often include unnatural, inconsistent, or unreadable text or logos.

(10) **Physical Laws:** Check whether the content of the image violates physical laws.

(11) **Faces:** Check whether the image content includes faces; AI-generated images often exhibit unnatural facial structures. This includes but is not limited to: 1. Overly perfect facial features lacking the imperfections of real faces. 2. Blurred and distorted edges and facial features. AI-generated faces may appear overly rough, leading to blurriness and distortion at the facial edges. Additionally, you can pay attention to the facial features generated by AI, which often have blurry traces and exaggerated, distorted expressions that differ from real-world faces. 3. Exaggerated and unnatural facial expressions. 5. Inconsistent skin color and texture. 6. Makeup. Overly perfect makeup or makeup inconsistent with age and gender are common flaws in AI-generated images. 7. Body accessories and clothing.

(12) **Body Structure:** Check whether the image content includes body structures; generally, AI-generated images will include some anatomical errors. The body parts you need to consider include head, neck, body, arm, hand, leg, and foot.

(13) **Common Sense:** Check whether there are any violations of common sense in the image.

The following four options are some annotations of this image:

- A. {text1}
- B. {text2}
- C. {text3}
- D. {text4}
- E. {text5}
- F. {text6}
-

The ground truth text you can refer to is as follows: {Ground Truth Text}

Please examine the provided image attentively and serve as an unbiased judge in assessing the quality of the response from an AI assistant regarding the instruction. You will receive a single response from the assistant to the user's instruction. The following are some notices and scoring criteria:

Notices:

1. Your assessment should identify whether the assistant effectively adheres to the user's instructions and addresses the user's inquiry.
2. In your evaluation, weigh factors such as relevance, accuracy, comprehensiveness, creativity, and the granularity of the responses.
3. Do not allow the length of the responses to influence your evaluation.
4. Do not favor certain names or positions of the assistants. Be as objective as possible.
5. If the available ground truth text is not empty, please refer to the ground truth text. However, the ground truth text may not be entirely correct, so you need to make a judgment and then score the different responses.

Criteria: Use scores to show the quality of the response. Here is the detailed scoring rubric for evaluating the quality of responses from AI assistants:

Poor (1): The response significantly deviates from the user's instruction and fails to address the query effectively. It shows a lack of relevance, accuracy, and comprehensiveness. Creativity and granularity are absent or poorly executed.

Fair (2): The response addresses the user's instruction partially, with evident shortcomings in relevance, accuracy, or comprehensiveness. It lacks depth in creativity and granularity, indicating a superficial understanding of the user's inquiry.

Average (3): The response adequately addresses the user's instruction, showing a fair level of relevance, accuracy, and comprehensiveness. It reflects a basic level of creativity and granularity but may lack sophistication or depth in fully capturing the user's inquiry.

Good (4): The response is well-aligned with the user's instruction, demonstrating a high degree of relevance, accuracy, and comprehensiveness. It shows creativity and a nuanced understanding of the topic, with a detailed granularity that enhances the response quality.


Excellent (5): The response perfectly adheres to the user's instruction, excelling in relevance, accuracy, comprehensiveness, creativity, and granularity. It provides an insightful, detailed, and thorough answer, indicating a deep and nuanced understanding of the user's inquiry.

Please score according to the following format:

- A. Score
- B. Score
- C. Score
- D. Score
- E. Score
- F. Score
-




Figure 8. Prompt for cross-model validation and evaluation using state-of-the-art multimodal large language models.



[Face]

Expert Images



You are an AI visual assistant that can help humans analyze images that may have been generated by AI. You will receive an AI-generated image. Your task is to uncover clues that suggest the image may be AI-generated. You need to perform the following two tasks:

1. Please describe the image in detail. At the beginning of the description, state, "This is a fake image."
2. Identify clues that differentiate this image from real images and describe them. You can consider the following angles, but are not limited to them, to find signs that the image is AI-generated.

Check whether the image content includes faces; AI-generated images often exhibit unnatural facial structures. This includes but is not limited to:


1. Overly perfect facial features lacking the imperfections of real faces.
2. Blurred and distorted edges and facial features. AI-generated faces may appear overly rough, leading to blurriness and distortion at the facial edges. Additionally, you can pay attention to the facial features generated by AI, which often have blurry traces and exaggerated, distorted expressions that differ from real-world faces.
3. Exaggerated and unnatural facial expressions.
5. Inconsistent skin color and texture.
6. Makeup. Overly perfect makeup or makeup inconsistent with age and gender are common flaws in AI-generated images.
7. Body accessories and clothing. Please note that you do not need to consider every point mentioned above in your response.

You are also encouraged to look for other clues related to AI-generated faces that were not mentioned. You only need to describe the clues you are most confident about. Please select the most critical clues that indicate this is an AI-generated image and describe them. You can respond in the following format:

To determine that this image is AI-generated, let's examine a few critical clues: 1.xxxx 2.xxxx 3.xxxx ...


Please provide the most relevant clues that indicate this is an AI-generated image, avoiding vague or uninformative responses.

Figure 9. Specialist Prompt for AI-generated images containing face defects.



[Body]

Expert Images



This is an AI-generated image.

The area enclosed in **the red box** contains anatomical errors: **{Abnormal Types}**.

You can consider aspects such as abnormal numbers, distortion, and proportion anomalies. Please focus solely on the errors within the area enclosed in the red box and describe them in relation to the content of the image.

You can respond in the following format:

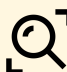
The anatomical errors present in this image are as follows:

1. **Abnormal Type1:** xxxx
2. **Abnormal Type2:** xxxx
3. **Abnormal Type3:** xxxx

...


Please note that the term 'red box' should not appear in your response. Please provide the most relevant clues that indicate this is an AI-generated image, avoiding vague or uninformative responses.

Figure 10. Specialist Prompt for AI-generated images containing body defects.



[Text&Logos]

Expert Images



This image contains areas with text, symbols, or logos.

Since this is an AI-generated image, there may be some issues with them, such as **blurriness, distortion, or meaningless characters**.

Please check all areas of the image that contain **similar text, symbols, or logos**, and identify some areas with relatively minor issues that you are confident about, while categorizing the rest as problematic. Additionally, describe the issues for each problematic area of text, symbols, or logos (such as **blurriness, distortion, or meaningless characters, etc**).

You can respond in the following format:

"In the provided image, there are several areas containing text and symbols.
Let's categorize them and identify any issues accordingly:

Text: xxx Issue: xxx
Text: xxx Issue: xxx "

For the problematic parts, please express them as traces related to AI generation, and do not attempt to forcefully identify text in those areas.

Figure 11. Specialist Prompt for AI-generated images containing defects in text&logos.

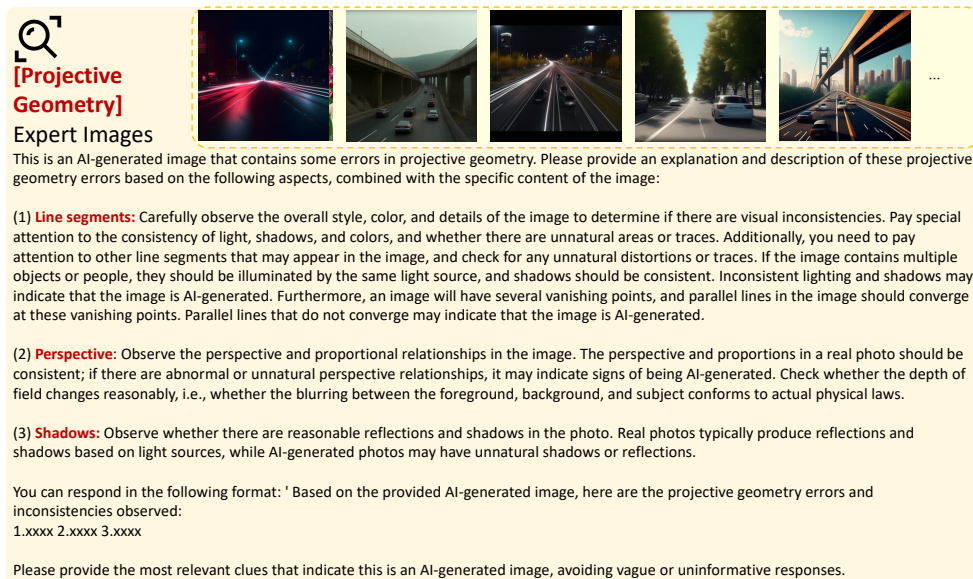


Figure 12. Specialist Prompt for AI-generated images containing defects in projective geometry.

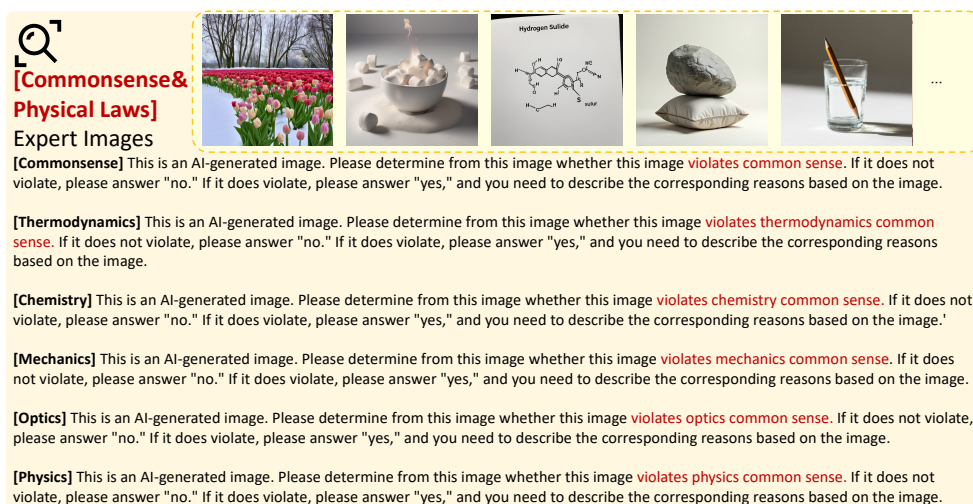


Figure 13. Specialist prompt for AI-generated images containing defects in commonsense and physical laws.

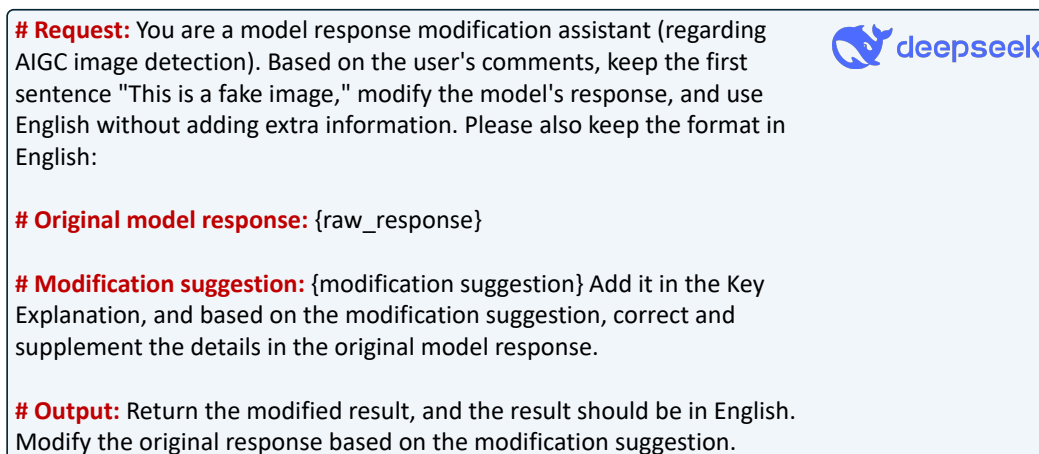


Figure 14. The original response of the model and the revision suggestions are input into the query prompt of DeepseekV3 [55].



Figure 15. Qualitative Results of the Test Set for \mathcal{P}_3 .

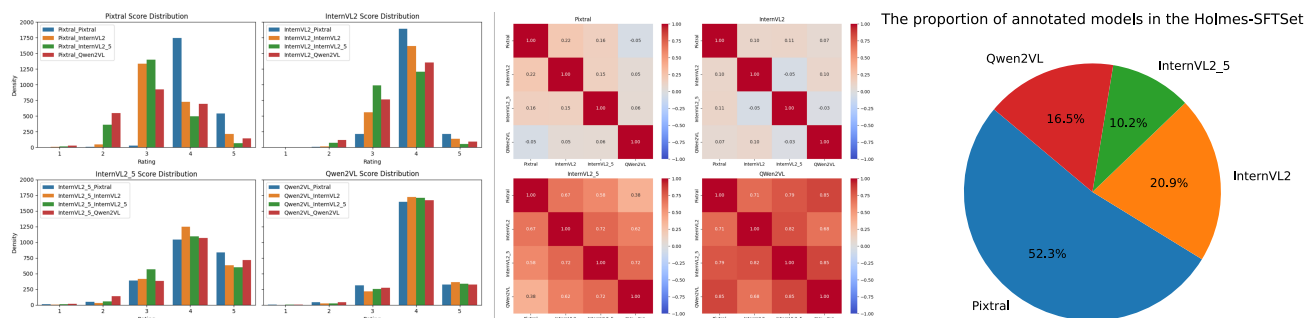


Figure 16. The score distribution of Multi-Expert Jury ratings, the correlation heatmap of the ratings, and the proportion of annotations by each MLLM Expert in Holmes-SFTSet.

You need to modify the following text and make some suggestions until you are basically satisfied. We will use mlmm to modify our explanation based on your suggestions. You can refer to the document: [click to visit](#)



☐ This question has a problem and cannot be answered

Source title ID

Topic ID

Labeler

You need to modify the following text and make some suggestions until you are basically satisfied. We will use mlmm to modify our explanation based on your suggestions. You can refer to the document: [click to visit](#)

Based on the provided fake image, here are the reasons why this image is fake:

- Line segments** : The shadows and light distribution on the clock faces and the surrounding architecture are inconsistent. The light source seems to come from different directions, which is physically impossible in a real-world scenario. Additionally, the parallel lines of the wooden beams do not converge correctly at a vanishing point, indicating a distortion typical of AI-generated images.
- Edges** : There are noticeable discontinuities and inconsistencies in the edges of the clock hands and the Roman numerals. The edges appear overly sharp in some areas and blurred in others, which is a common artifact in AI-generated images.
- Texture** : The texture of the wooden beams and the clock faces appears unnaturally smooth and lacks the detailed grain and imperfections that would be present in a real photograph. The repetitive pattern of the wooden beams also suggests an AI-generated texture.
- Distortion** : The clock faces exhibit slight distortions, particularly around the edges, which are not consistent with the perspective and proportions of a real clock tower. The hands of the clock also appear slightly warped, which is unnatural.
- Overall Hue** : The overall tone of the image is overly vibrant and lacks the natural variations in lighting and color that would be present in a real photograph. The contrast between the dark and light areas is too stark, indicating artificial enhancement.
- Clarity** : The image shows signs of unnatural pixel blurriness, particularly around the edges of the clock faces and the hands. This suggests that the image may have been generated by AI and then compressed, leading to artifacts.
- Perspective** : The perspective of the image is slightly off, with the clock faces appearing disproportionately large compared to the surrounding architecture. The depth of field does not change reasonably, with both the foreground and background appearing equally sharp, which is unnatural.
- Shadows** : The shadows cast by the clock hands and the Roman numerals do not correspond to a single light source. The shadows are inconsistent and do not follow the natural laws of light and shadow, indicating that the image is AI-generated.

By carefully analyzing these aspects, it becomes evident that the image is likely AI-generated.

You can modify the above explanation of this picture. You do not need to modify the part of the authenticity judgment because it is accurate. You can do it from the following perspectives: 1. Based on your observations and explanations, what kind of explanation do you want the model to have? If it is not in the explanation, please add it. 2. Complete your additional explanation suggestions and the modification suggestions for the original explanation from top to bottom in the order of your preference.

☐ This question has a problem and cannot be answered

Source title ID

Topic ID

Labeler

Figure 17. The annotation interface for obtaining preference samples with modifications suggested by human experts.

You need to compare the responses of option A and option B according to the requirements of the question. These two responses are generated by different models and contain some descriptions about the authenticity of the pictures. You need to evaluate objectively.



A.

Based on the provided fake image, here are the reasons why this image is fake:

1. **Texture:** The image exhibits unnatural textures, especially in the clothing and hair of the people. The textures appear blurry and lack detail, which is common in AI-generated images.
2. **Edges:** The edges of the people and the food truck are not as sharp as they would be in a real image. There are noticeable pixelation and blurriness around the edges of the objects.

☐ This question has a problem and cannot be answered

Source ID Title ID

You need to compare the responses of option A and option B according to the requirements of the question. These two responses are generated by different models and contain some descriptions about the authenticity of the pictures. You need to evaluate objectively.

4. **Distortion:** The proportions of some individuals seem distorted. For instance, the person on the right seems to have an awkwardly large head compared to the body, which is not typical in real images.
5. **Overall Hue:** The colors in the image are overly vibrant and have an unnatural saturation level. The blue of the sky and the yellow of the food truck are particularly striking and do not blend naturally with the rest of the image.
6. **Clarity:** The image has a slight blurriness around key areas like the people and the food truck, which is unusual in a high-quality photograph.
7. **Perspective:** The perspective of the image is slightly off. The people and the food truck do not align properly in terms of perspective, giving an unrealistic spatial relationship.
8. **Shadows:** There are inconsistencies in the shadows cast by the people and the food truck. The shadows do not match the direction and intensity of the light source, indicating artificial manipulation.
9. **Text:** The text on the food truck is blurry and difficult to read, which is unusual for a clear and well-defined image. This blurriness suggests it may have been digitally added.
10. **Physical Laws:** The image violates physical laws in terms of lighting and shadow consistency. The light source direction does not correspond with the shadows cast by the people and objects.
11. **Faces:** The faces of the people in the group appear to have unnatural features, with overly smooth skin and inconsistent lighting on their faces, indicating they may have been generated or manipulated by AI.
12. **Body Structure:** Some individuals in the image exhibit anatomical errors, such as disproportionate body parts. For instance, the person on the left has an awkwardly positioned arm, which looks unnatural.
13. **Common Sense:** The overall composition of the image does not make logical sense. The people are clustered in an unnatural manner, and the food truck seems out of place in the setting.

These clues collectively indicate that the image is AI-generated.

Your evaluation should determine whether the assistant effectively follows the user's instructions and addresses the user's inquiry. Do not let the length of the response influence your evaluation. Do not favor certain assistants by name or position. Be as objective as possible. In your evaluation, weigh the following factors: relevance, accuracy, comprehensiveness, creativity, and granularity of the response: **Relevance:** The judge's decisions correspond directly to the instructions or standards provided. Each judgment is relevant to the case at hand and does not stray into irrelevant areas. **Accuracy:** The judge's decisions consistently conform to the established rules or guidelines. These guidelines are clearly understood and correctly applied in each judgment.

* **Comprehensiveness:** The judge considers all necessary aspects and evidence related to each case. The judge covers every relevant point in the guidelines in his or her evaluation. **Creativity:** The judge demonstrates the ability to think creatively when dealing with complex or ambiguous situations. This includes providing insightful, constructive feedback or solutions, even if these are not explicitly covered in the guidelines. **Granularity of the response:** The judge provides detailed and specific reasons for each decision. This includes a thorough analysis of how each aspect of the guidelines applies to the current case or situation. Please sort your two options A and B. If A's answer is better, choose A. If B's answer is better, choose B. If they are tied, choose C.

☐ A: A is better ☒ B: B is better ☐ C: those

☐ This question has a problem and cannot be answered

Source ID Title ID

Figure 18. The interface for evaluating sample selection of human preferences in Arenas.

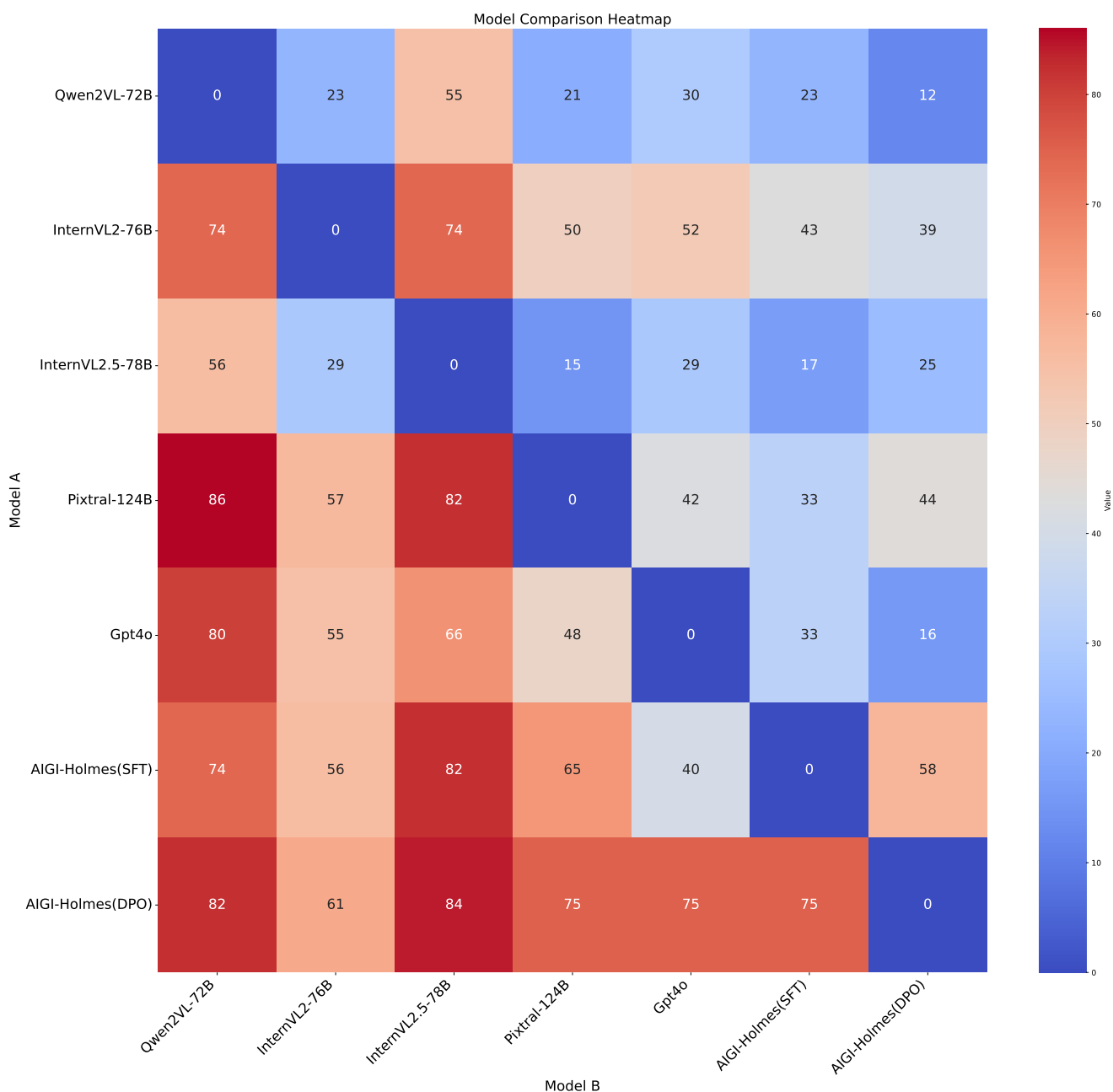


Figure 19. A heatmap of the winning counts in pairwise scoring among state-of-the-art multimodal language models (MLLMs)

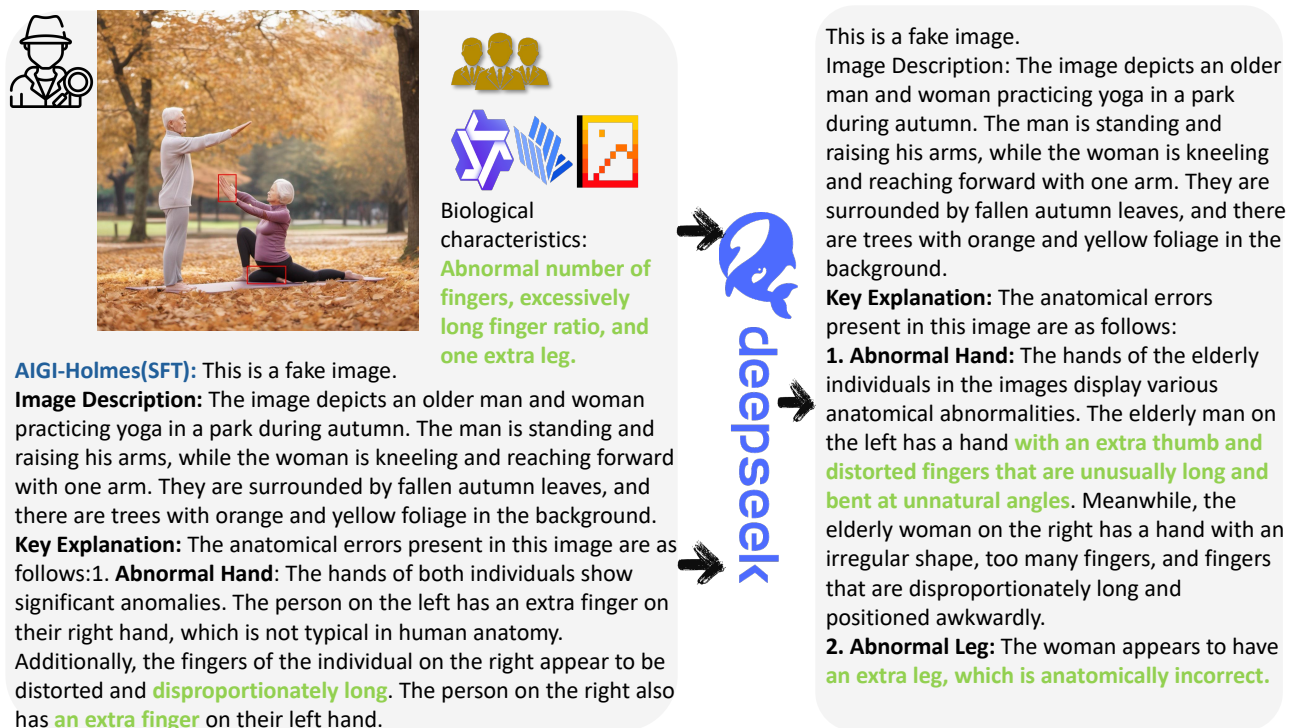


Figure 20. The schematic flowchart of the DeepseekV3 [55] generating human-aligned preference samples based on modification suggestions and original responses.



QwenVL-72B: 3.75



This is a fake image.

Image Description: The image depicts a sheep sitting on the grass in a field. The sheep appears to be in an unusual and unnatural position, with its body contorted in a way that is not typical for a sheep. The background shows a few other sheep lying on the ground, and there are trees in the distance. The overall scene is serene, but the sheep's posture raises questions about the image's authenticity.

Based on the provided fake image, here are the reasons why this image is fake:

1. **Body Structure:** The sheep's body structure is highly unusual. The way its body is contorted and supported by its legs does not conform to the natural anatomy of a sheep. This unnatural posture is a strong indicator that the image is AI-generated.
2. **Texture:** The texture of the sheep's wool appears overly smooth and uniform, lacking the natural variations and imperfections typically found in real fur or wool. This homogeneity in texture is a common trait in AI-generated images.
3. **Distortion:** The overall scene, particularly the sheep's body, exhibits signs of distortion. The proportions and shapes do not align with real-world physics and anatomy, suggesting that the image has been manipulated or generated by AI.

Figure 21. Annotated sample and its score in Holmes-SFTSet for Qwen2VL-72B.



InternVL2-76B: 4.00



This is a fake image.

Image Description: The image shows two women with dark hair and light skin standing closely together against a green background. They are smiling and looking directly at the camera. One woman is wearing a light blue shirt, and the other is wearing a beige top. The background is solid green, and there is no other significant detail present.

Based on the provided fake image, here are the reasons why this image is fake:

- 1. Line Segments:** There are some inconsistencies in the lighting and shadows. The lighting on the faces of the two women is **uniform and lacks natural variation**, suggesting that the light source is not consistent with a real-world setting. The shadows under their chins and around their necks are slightly blurry and unnatural. The edges around their hair and faces appear **slightly too smooth and lack the natural texture** typically seen in real photos.
- 2. Edges:** The edges around the women's hair and shoulders blend unnaturally with the background, and there are areas where the edges appear too smooth, indicating potential pixel manipulation. The outlines of their faces and hair seem unnaturally clean and lack the normal imperfections you would expect in a real photograph.
- 3. Texture:** **The texture of their clothing lacks the natural folds and creases you would see in real fabric, especially noticeable in the blue shirt of the woman on the right. The fabric appears too flat and smooth**, hinting at AI generation.
- 4. Distortion:** There is a slight distortion in the alignment of the faces, particularly around the eyes and mouth. The facial features, while detailed, have a hint of blurriness that suggests they may have been generated or heavily edited.
- 5. Overall Hue:** The overall hue of the image is slightly too vibrant, which is often a characteristic of AI-generated images. The colors are overly saturated and appear unnatural, especially in the green background.
- 6. Clarity:** The image lacks fine details that would be present in a high-resolution photograph, such as individual strands of hair or subtle skin textures. **The quality of the image is slightly blurred, particularly around the edges**, indicating potential compression artifacts.
- 7. Perspective:** The perspective is consistent and natural, **but the proportions of the women's faces are slightly off, with their facial features appearing exaggerated and too perfect**, which is a common trait of AI-generated images.
- 8. Shadows:** The shadows under their chins are inconsistent and lack depth, suggesting they were artificially added or edited rather than naturally occurring.
- 9. Text:** There is no text in this image, so this point is not applicable.
- 10. Physical Laws:** There are no apparent violations of physical laws in this image, as the image is relatively simple and focused on the two women.
- 11. Faces:** Their facial features are too perfect and symmetrical, lacking natural imperfections. The skin texture is unnaturally smooth, and the makeup appears overly perfect and uniform.
- 12. Body Structure:** The body structures appear generally correct, but the lack of natural imperfections in their skin and the smoothness of their features suggest AI generation.
- 13. Common Sense:** There is nothing inherently illogical or nonsensical in the image, but the overall perfection and smoothness of the image are not consistent with typical photographs.

Figure 22. Annotated sample and its score in Holmes-SFTSet for InternVL2-76B.



InternVL2.5-78B: 4.00



This is a fake image.

Image Description: The image depicts a yellow bus at a bus stop with several people boarding or waiting to board. The bus has a sign indicating its destination, and the background shows a cityscape with buildings and trees. The sky is partially cloudy, and the overall scene appears to be a typical urban setting.

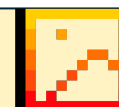
Based on the provided fake image, here are the reasons why this image is fake:

1. **Line segments:** The shadows cast by the people and the bus do not align consistently with the light source. The shadows appear to be inconsistent and unnatural, suggesting that they may have been added or altered.
2. **Edges:** The edges of the bus and some of the people appear slightly blurred and lack the sharpness typically found in real photographs. This suggests that these elements may have been generated or added separately.
3. **Texture:** The texture of the road and the bus's surface appears overly smooth and lacks the natural variations found in real images. This can be a sign of AI generation.
4. **Distortion:** The proportions of the bus and the people seem slightly off, with the bus appearing somewhat elongated and the people's sizes varying unnaturally.
5. **Overall Hue:** The overall tone of the image is slightly oversaturated, particularly in the yellow of the bus and the blue of the sky, which can be characteristic of AI-generated images.
6. **Clarity:** The image lacks the natural grain and noise that would be present in a real photograph, especially in the sky and the distant buildings.
7. **Perspective:** The perspective of the bus and the people boarding it does not seem to follow realistic depth cues, with some elements appearing flat or improperly scaled.
8. **Shadows:** The shadows on the ground do not match the positions of the people or the bus, indicating that they may have been added separately or improperly rendered.
9. **Text:** The text on the bus sign is not entirely clear and appears to be a mix of real and generated elements, which is common in AI-generated images.
10. **Physical Laws:** The way the light interacts with the surfaces in the image, particularly on the bus and the people, does not follow real-world physics, suggesting that the lighting may have been artificially manipulated.
11. **Faces:** The faces of the people are difficult to make out and appear somewhat blurred, which is a common issue in AI-generated images where facial features are not well-rendered.
12. **Body Structure:** The body structures of the people boarding the bus are not clearly defined and appear somewhat distorted, indicating potential issues with AI generation of human forms.

Figure 23. Annotated sample and its score in Holmes-SFTSet for InternVL2.5-78B.



Pixtral-124B: 4.75



This is a fake image.

Image Description: The image shows a tennis racket resting on a tennis court with a tennis ball near the strings. The tennis ball has the word "P-E-I-A-N" printed on it in blue letters. The court surface appears to be a hard court, and the racket's handle is visible in the foreground.

Based on the provided fake image, here are the reasons why this image is fake:

1. **Text:** The text "P-E-I-A-N" on the tennis ball is not a recognizable brand or logo commonly seen on tennis balls. Real tennis balls typically have brand names like "Wilson," "Penn," or "Dunlop." The presence of an unrecognizable or nonsensical text is a strong indicator that the image is AI-generated.
2. **Texture:** The texture of the tennis ball appears overly smooth and lacks the typical fuzziness associated with real tennis balls. This unnatural texture is a common artifact in AI-generated images.
3. **Shadows and Lighting:** The shadows cast by the racket and the ball do not seem consistent with a single light source. The shadow of the racket appears to be coming from a different direction than the shadow of the ball, which is physically impossible in a real photograph.
4. **Edges:** The edges of the racket strings and the tennis ball show slight unnatural pixelation and blurriness, which are often signs of AI generation. Real photographs typically have sharper and more defined edges.
5. **Overall Hue:** The overall color tone of the image appears somewhat dull and lacks the vibrancy typically seen in real photographs taken outdoors. This inconsistency in hue is another clue that the image may be AI-generated.

These clues collectively suggest that the image is likely AI-generated rather than a real photograph.

Figure 24. Annotated sample and its score in Holmes-SFTSet for Pixtral-124B.



Chose

This is a fake image.

Image Description: The image depicts an electric car charging station with multiple charging points. There are three charging stations visible, with two white cars parked and being charged. The environment appears to be an outdoor parking area with trees and a building in the background. The charging cables are connected to the cars, and one cable is lying on the ground.

Based on the provided fake image, here are the reasons why this image is fake:

1. **Line segments:** The shadows cast by the cars and the charging stations are not consistent with the direction of the light source. The shadows appear to be going in different directions, which is physically impossible in a real-world setting. Additionally, the parallel lines of the parking spaces do not converge correctly, indicating a possible AI-generated distortion.
2. **Edges:** The edges of the charging stations and the cars appear overly smooth and lack the natural irregularities that are typically present in real photographs. This smoothness suggests AI generation.
3. **Texture:** The texture of the ground and the trees in the background appear unnaturally blurry and lack the detail that would be expected in a real photo. Additionally, the repetitive pattern of the leaves on the trees suggests AI-generated textures.
4. **Distortion:** The shape of the cars and the charging stations exhibit slight distortions that are not typical of real objects. For example, the wheels of the cars appear slightly elongated, which is not physically accurate.
5. **Overall Hue:** The overall tone of the image is overly vibrant, with colors that appear too saturated and unnatural. This is a common characteristic of AI-generated images.
6. **Clarity:** The image shows signs of unnatural pixel blurriness, especially around the edges of the cars and the charging stations. This lack of clarity is often seen in AI-generated images.
7. **Perspective:** The perspective of the image is inconsistent, with the depth of field not changing reasonably. The foreground and background do not blur naturally, indicating a possible AI generation.
8. **Shadows:** The shadows cast by the objects in the image are not consistent with the light source. The shadows appear to be going in different directions, which is not possible in a real-world setting.
9. **Physical Laws:** The charging cables appear to be too flexible and do not follow the natural laws of physics. They are lying on the ground in an unnatural manner, which is not typical of real charging cables.
10. **Common Sense:** The placement of the charging stations and the cars is not practical. In a real-world setting, the charging stations would be placed in a more organized manner, and the cars would be parked more neatly. The current arrangement does not make sense in a practical setting.

These clues suggest that the image is likely AI-generated.

Reject

This is a real image.

Image Description: This image depicts an electric car parked in a designated parking spot, with multiple electric vehicle charging stations visible. The charging stations are installed next to the parking spots, with cables extending from the stations to the vehicles. The scene appears to be outdoors, with trees and a building visible in the background. The car is connected to one of the charging cables, indicating that it is being charged.

Based on the provided real image, here are the reasons why this image is real:

1. **Line Segments:** The parallel lines of the charging stations and the parking lines converge naturally toward a vanishing point, indicating a realistic perspective. The lighting and shadows are consistent, with the shadows of the charging stations and the car aligning with a single light source, likely the sun. There are no visible unnatural distortions or inconsistencies in the lines or shadows.
2. **Edges:** The edges of the car, charging stations, and other objects in the image are clear and continuous, without any discontinuous or inconsistent edges that might suggest AI generation. The pixel distribution around the edges is natural and smooth.
3. **Texture:** The textures of the car, charging stations, and surrounding environment are realistic and detailed. There are no blurry or unclear textures, nor are there repetitive texture patterns that would indicate AI generation. The surfaces of the objects appear natural and consistent with real-world materials.
4. **Distortion:** The objects in the image, including the car and charging stations, exhibit no unnatural distortions. The shapes and proportions of the objects conform to real-world expectations, with no anomalies or unnatural warping.
5. **Overall Hue:** The overall tone of the image is balanced and consistent, with no areas that are overly vibrant or dull. The colors are natural and reflect the expected hues of an outdoor scene with electric vehicles and charging infrastructure.
6. **Clarity:** The image has a high resolution with no visible compression artifacts. There is no unnatural pixel blurriness or jagged edges, indicating that the image is not AI-generated.
7. **Perspective:** The perspective and proportional relationships in the image are consistent. The depth of field changes reasonably, with the foreground, background, and subject all in focus, conforming to actual physical laws.
8. **Shadows:** The reflections and shadows in the image are reasonable and consistent with the light source. The shadows cast by the car and charging stations are natural and align with the direction of the sunlight.
9. **Text:** There is no visible text or logos in the image that could be analyzed for inconsistencies.
10. **Physical Laws:** The content of the image adheres to physical laws, with no violations of natural principles.
11. **Faces:** The image does not contain any faces, so this criterion is not applicable.
12. **Body Structure:** The image does not contain any human bodies, so this criterion is not applicable.
13. **Common Sense:** The scene depicted in the image is realistic and adheres to common sense, with no anomalies or inconsistencies that would suggest AI generation. The presence of electric vehicles and charging stations in an outdoor setting is plausible and aligns with real-world scenarios.

electric vehicles and charging stations in an outdoor setting is plausible and aligns with real-world scenarios.

Figure 25. The presentation of preference sample pairs in \mathcal{D}_1 .



Chose

This is a fake image.

Image Description: The image depicts a teddy bear sitting upright with a metallic rectangular block balanced precariously on its head. The teddy bear has a light brown, soft-looking texture, and the metallic block appears to be made of aluminum foil or a similar material, with visible creases and folds. The background is a plain, light-colored surface, and the lighting is soft and diffused, casting subtle shadows.

Key Explanation: This image violates common sense in terms of mechanics and physics. Here are the reasons:

1. **Balance and Stability:** The metallic block is precariously balanced on the head of the teddy bear. Given the apparent shape and texture of the teddy bear, it does not provide a stable or flat surface for the block to rest on. The block would likely slip off due to the soft and uneven surface of the teddy bear.
2. **Weight Distribution:** The metallic block appears to be quite heavy, especially if it is made of metal. The teddy bear, being soft and plush, would not be able to support the weight of the block without collapsing or deforming significantly.
3. **Material Properties:** The teddy bear is made of soft material, likely foam or fabric. Such materials are not designed to bear the weight of heavy, rigid objects like a metallic block without being damaged.

Based on the provided fake image, here are the reasons why this image is fake:

1. **Physical Laws:** The most obvious clue is the violation of physical laws. The metallic block is precariously balanced on the teddy bear's head, which is highly improbable in real life. The block appears to be too heavy and unwieldy to be stable on such a soft and uneven surface. This defies the laws of physics and common sense, indicating that the image is likely AI-generated.
2. **Shadows:** The shadows in the image do not fully align with the expected lighting. The shadow cast by the metallic block on the teddy bear is not as pronounced as it should be, given the apparent weight and size of the block. Additionally, the shadow of the teddy bear on the surface is very faint, suggesting inconsistent lighting.
3. **Texture:** The texture of the teddy bear appears overly uniform and lacks the natural variations and imperfections that would be present in a real plush toy. The fur appears too perfect and evenly distributed, which is a common characteristic of AI-generated images.
4. **Perspective:** The perspective of the image seems slightly off. The teddy bear and the metallic block are positioned in a way that suggests a top-down view, but the background and the surface they are on do not align perfectly with this perspective, creating a slight disconnect.
5. **Common Sense:** From a common-sense standpoint, the scenario depicted is highly unlikely. It is rare to see a teddy bear supporting a heavy metallic block in such a manner, especially in a stable and balanced way. This unusual and impractical setup is a strong indicator that the image is artificial.

These clues collectively suggest that the image is AI-generated rather than a real photograph.

Reject

This is a fake image.

Image Description: The image depicts a teddy bear sitting upright with a metallic rectangular block balanced precariously on its head. The teddy bear has a light brown, soft-looking texture, and the metallic block appears to be made of aluminum foil or a similar material, with visible creases and folds. The background is a plain, light-colored surface, and the lighting is soft and diffused, casting subtle shadows.

Key Explanation: This is a fake image. Here are the reasons:

1. **Physical Laws:** The most obvious clue is the violation of physical laws. The metallic block is precariously balanced on the teddy bear's head, which is highly improbable in real life. The block appears to be too heavy and unwieldy to be stable on such a soft and uneven surface. Additionally, there is no visible deformation of the teddy bear under the weight of the block, which defies natural laws of physics and common sense, indicating that the image is likely AI-generated.
2. **Shadows:** The shadows in the image do not fully align with the expected lighting. The shadow cast by the metallic block on the teddy bear is not as pronounced as it should be, given the apparent weight and size of the block. Additionally, the shadow of the teddy bear on the surface is very faint, suggesting inconsistent lighting.
3. **Texture:** The texture of the teddy bear appears overly uniform and lacks the natural variations and imperfections that would be present in a real plush toy. The fur appears too perfect and evenly distributed, which is a common characteristic of AI-generated images. Moreover, the texture is unnaturally clear, further indicating artificial creation.
4. **Perspective:** The perspective of the image seems slightly off. The teddy bear and the metallic block are positioned in a way that suggests a top-down view, but the background and the surface they are on do not align perfectly with this perspective, creating a slight disconnect.
5. **Common Sense:** From a common-sense standpoint, the scenario depicted is highly unlikely. It is rare to see a teddy bear supporting a heavy metallic block in such a manner, especially in a stable and balanced way. This unusual and impractical setup is a strong indicator that the image is artificial.

These clues collectively suggest that the image is AI-generated rather than a real photograph.

Figure 26. The presentation of preference sample pairs in \mathcal{D}_2 .



Qwen2VL-72B:

The image is fake due to the following reasons:

1. Line segments and edges: Unnatural distortions and inconsistencies, especially around the laptops and hands, with blurred edges common in AI-generated images.
2. Texture: Overly smooth clothing and skin, lacking natural variations typical of real images.
3. Distortion: Warped laptops and inconsistently shaped hands, indicating unnatural distortions.
4. Overall hue: Overly vibrant tones with unnatural color variations, a common trait of AI-generated images.
5. Clarity: Unnatural pixel blurriness around edges, particularly on laptops and hands.
6. Perspective: Incorrect depth of field and proportional relationships, failing to align with real-world physics.
7. Shadows: Unnatural reflections and misaligned shadows, inconsistent with the light source.
8. Faces: Overly perfect facial features with blurred and distorted edges, typical of AI-generated images.
9. Body structure: Anatomical errors, such as distorted hands and misaligned feet.



AIGI-Holmes(SFT):

Reasons why the image is fake:

1. Texture: The couch and floor are unnaturally smooth, lacking real-life imperfections like fabric variations or floor scratches.
2. Distortion: Proportions are off; laptops appear too large for the boys' laps, and their bodies show slight elongation, especially in arms and legs.
3. Overall hue: Colors are overly vibrant and uniform, lacking natural variations typical in real photos.
4. Faces: Facial features are overly perfect, with smooth skin and slightly blurred edges, common in AI-generated faces.
5. Body structure: Anatomical errors, such as an unnaturally long neck on the left boy and disproportionately large hands on the right boy.
6. Common sense: Unusual details like both boys being barefoot on a couch and laptops placed unsupported on their laps, which is impractical.

These inconsistencies confirm the image is AI-generated.



AIGI-Holmes(DPO):

Key signs of it being AI-generated include:

1. Facial Features: Overly symmetrical and perfect, lacking natural imperfections.
2. Hands and Fingers: Distorted and lacking natural detail, particularly the boy on the right.
3. Clothing: Overly crisp and perfect designs, missing natural wrinkles.
4. Background: Artificial appearance with uniform lighting and shadows.
5. Shadows: Inconsistent and misaligned, suggesting multiple light sources.
6. Edges: Discontinuities, with some areas overly sharp and others blurred.
7. Texture: Unnaturally smooth, lacking detailed fabric or surface textures.
8. Distortion: Slight distortions in hand positions and laptop shapes.
9. Anatomy: Errors in body proportions and unnatural positioning of hands and feet.

These inconsistencies confirm the image is AI-generated.

Figure 27. Comparison of qualitative results between Qwen2VL-72B and our method.



InternVL-76B:

This is a fake image. Here are the reasons why:

1. Line segments: The lighting on the refrigerator and surrounding objects is inconsistent, especially the reflection on the refrigerator door and shadows under the shelves.
2. Edges: The edges of the refrigerator and countertop are unnatural and slightly blurred, particularly around the handles and edges.
3. Texture: The countertop and wood grain of the cabinets appear repetitive and artificial, lacking subtle variations.
4. Distortion: Objects on the countertop, like the circular item with handwritten text, show slight distortion, indicating they might have been superimposed.
5. Overall Hue: The image is too vibrant and lacks natural color variation, with overly crisp colors.
6. Clarity: There are areas of unnatural pixel blurriness, especially where different elements meet, and the overall clarity is inconsistent with a high-resolution photograph.
7. Perspective: The perspective is correct, but the depth of field seems unnatural, suggesting digital assembly.
8. Shadows: Shadows and reflections are generally consistent, but the reflection on the refrigerator's surface does not align with the light source and surrounding objects.



AIGI-Holmes(SFT):

1. Line segments: Shadows and lighting are inconsistent; shadows don't align with the light source, and refrigerator lines don't converge naturally.
2. Edges: Edges of objects are blurred and pixelated, unlike sharp edges in real photos.
3. Texture: Surfaces lack natural grain and imperfections, appearing overly smooth.
4. Distortion: Decorative item on the countertop is slightly distorted, a common AI flaw.
5. Overall Hue: Colors are too uniform and vibrant, lacking natural lighting variations.
6. Clarity: Unnatural pixel blurriness, especially around edges, is present.
7. Perspective: Refrigerator appears disproportionately large, showing perspective inconsistency.
8. Shadows: Shadows don't match the light source, and metallic surfaces lack reflections.

These clues suggest the image is AI-generated.



AIGI-Holmes(DPO):

Key Explanation: The provided image contains several areas with text and symbols that exhibit common AI-generated issues such as blurriness, distortion, and meaningless characters.

1. Text on the kitchen counter (top right corner): Blurry handwritten note, characters not clearly identifiable.
2. Symbols/Logos on the refrigerator door: Partially visible but distorted, unrecognizable.
3. Text on the kitchen counter (bottom right corner): Blurry and not clearly readable.

Figure 28. Comparison of qualitative results between Intern2VL-76B and our method.



InternVL2.5-78B:

This image is fake for several reasons:

1. The line segments and edges, especially around the hands and suit borders, appear unnatural and distorted.
2. The texture of the suits is overly smooth and lacks natural variance, with an artificial clarity and lack of detail in fabric folds.
3. The proportions and perspective of the figures are off, particularly in the arm and hand regions, showing unnatural distortions.
4. Shadows are inconsistent with the light source, with some areas missing shadows and others having unnatural shadows.
5. The overall color tone is overly vibrant and uniform, which is not typical of real photographs.
6. The suits do not conform naturally to the body, defying physical laws.
7. The faces, while relatively natural, lack imperfections and have a slight blurriness or smoothness around the edges.
8. The body structure shows anatomical inconsistencies, such as unnatural positioning and shape of the hands and arms.

These clues indicate that the image is likely AI-generated.



AIGI-Holmes(SFT):

Reasons why the image is fake:

1. Texture: The couch and floor are unnaturally smooth, lacking real-life imperfections like fabric variations or floor scratches.
2. Distortion: Proportions are off; laptops appear too large for the boys' laps, and their bodies show slight elongation, especially in arms and legs.
3. Overall hue: Colors are overly vibrant and uniform, lacking natural variations typical in real photos.
4. Faces: Facial features are overly perfect, with smooth skin and slightly blurred edges, common in AI-generated faces.
5. Body structure: Anatomical errors, such as an unnaturally long neck on the left boy and disproportionately large hands on the right boy.
6. Common sense: Unusual details like both boys being barefoot on a couch and laptops placed unsupported on their laps, which is impractical.

These inconsistencies confirm the image is AI-generated.



AIGI-Holmes(DPO):

Key signs of it being AI-generated include:

1. Facial Features: Overly symmetrical and perfect, lacking natural imperfections.
2. Hands and Fingers: Distorted and lacking natural detail, particularly the boy on the right.
3. Clothing: Overly crisp and perfect designs, missing natural wrinkles.
4. Background: Artificial appearance with uniform lighting and shadows.
5. Shadows: Inconsistent and misaligned, suggesting multiple light sources.
6. Edges: Discontinuities, with some areas overly sharp and others blurred.
7. Texture: Unnaturally smooth, lacking detailed fabric or surface textures.
8. Distortion: Slight distortions in hand positions and laptop shapes.
9. Anatomy: Errors in body proportions and unnatural positioning of hands and feet.

These inconsistencies confirm the image is AI-generated.

Figure 29. Comparison of qualitative results between InternVL2.5-78B and our method.



Pixtral-124B:

Based on the provided fake image, here are the reasons why this image is fake:

1. Texture Inconsistencies: The texture of the books and the bookshelf appears overly smooth and lacks the natural irregularities and imperfections that real objects typically exhibit. This is a common characteristic of AI-generated images.
2. Perspective Issues: The perspective of the bookshelf and its contents appears slightly off. The shelves do not seem to recede naturally into the background, suggesting a distortion that is often seen in AI-generated images.
3. Physical Laws Violation: The dog-shaped shadow is a clear violation of physical laws, as shadows should accurately represent the shape of the object casting them. This unnatural shadow shape is a strong indicator of AI manipulation.
4. Overall Hue: The overall hue of the image is somewhat dull and lacks the vibrancy typically seen in real photographs. This could be a result of AI generation, which sometimes struggles with accurate color reproduction.
5. Edges and Lines: The edges of the bookshelf and the books within it appear slightly blurred and lack the sharpness expected in a real photograph. This blurriness is a common artifact in AI-generated images.

These clues collectively suggest that the image is AI-generated and not a real photograph.



AIGI-Holmes(SFT):

Based on the provided fake image, here are the reasons why this image is fake:

1. Shadows: The shadow of the stack of books falling off the shelf is unnatural. In a real-world scenario, the books would likely fall straight down rather than floating in the air as a single shadow. Additionally, the shadow's shape and angle do not accurately reflect the light source, suggesting it was artificially generated.
 2. Perspective: The perspective of the bookshelf and the falling books is slightly off. The books appear to be defying gravity, as they are not aligned with the expected direction of fall. This inconsistency in perspective is a common flaw in AI-generated images.
 3. Physical Laws: The way the books are stacked and the manner in which they are falling violate basic physical laws. The books appear to be floating and not following a natural trajectory, which is physically impossible.
 4. Clarity: The edges of the books and the bookshelf are overly sharp and clear, lacking the natural blur and imperfections that would be present in a real photograph. This hyper-realistic clarity is often a sign of AI generation.
- By considering these factors, it becomes evident that the image is likely AI-generated.



AIGI-Holmes(DPO):

Key Explanation: This image violates optics common sense. Here are the reasons:

1. Inconsistent Shadows: The shadows of the books and the shelf on the wall do not align correctly with the light source. In a real-world scenario, the shadows should be cast in a direction opposite to the light source. Here, the shadows appear to be going in multiple directions, which is not physically accurate.
2. Unnatural Shadow Intensity: The intensity and shape of the shadows are inconsistent. Some shadows are very dark while others are very faint, even though the light source should be uniform.
3. Lack of Perspective: The perspective of the shadows does not match the perspective of the objects. The shadows on the wall should be elongated and stretched out in a way that aligns with the angle of the light source, but they appear flat and two-dimensional.
4. Unrealistic Lighting: The lighting on the bookshelf and the wall does not match the direction and intensity of the shadows. The light seems to be coming from multiple directions, which is not possible with a single light source.

Figure 30. Comparison of qualitative results between Pixtral-124B and our method.