



AEGIS: A Holistic Benchmark for Evaluating Forensic Analysis of AI-Generated Academic Images

Bo Zhang* Tzu-Yen Ma* Zichen Tang Junpeng Ding Zirui Wang
 Yizhuo Zhao Peilin Gao Zijie Xi Zixin Ding Haiyang Sun Haocheng Gao
 Yuan Liu Liangjia Wang Yiling Huang Yujie Wang Yuyue Zhang
 Ronghui Xi Yuanze Li Jiacheng Liu Zhongjun Yang Haihong E[†]
 Beijing University of Posts and Telecommunications



bupt-reasoning-lab.github.io/AEGIS



BUPT-Reasoning-Lab/AEGIS



BUPT-Reasoning-Lab/AEGIS

Abstract

We introduce **AEGIS**, A holistic benchmark for Evaluating forensic analysis of AI-Generated academic ImageS. Compared to existing benchmarks, **AEGIS** features three key advances: (1) **Domain-Specific Complexity**: covering seven academic categories with 39 fine-grained subtypes, exposing intrinsic forensic difficulty, where even GPT-5.1 reaches 48.80% overall performance and expert models achieve only limited localization accuracy (IoU 30.09%); (2) **Diverse Forgery Simulations**: modeling four prevalent academic forgery strategies across 25 generative models, with 11 yielding average forensic accuracy below 50%, showing that forensics lag behind generative advances; and (3) **Multi-Dimensional Forensic Evaluation**: jointly assessing detection, reasoning, and localization, revealing complementary strengths between model families, with multimodal large language models (MLLMs) at 84.74% accuracy in textual artifact recognition and expert detectors peaking at 79.54% accuracy in binary authenticity detection. By evaluating 25 leading MLLMs, nine expert models, and one unified multimodal understanding and generation model, **AEGIS** serves as a diagnostic testbed exposing fundamental limitations in academic image forensics.

1 Introduction

With the rapid advancement of generative models, AI-generated images have reached high visual fidelity, raising concerns about the trustworthiness of visual evidence. In response, image forensics has evolved into multiple detection methods, including frequency-domain (Corvi et al., 2023), diffusion-process (Chen et al., 2024a), patch-level (Chen et al., 2024b), and alignment-based approaches (Sundara Rajan et al., 2025). Mean-

*Equal contribution.

[†]Corresponding author.

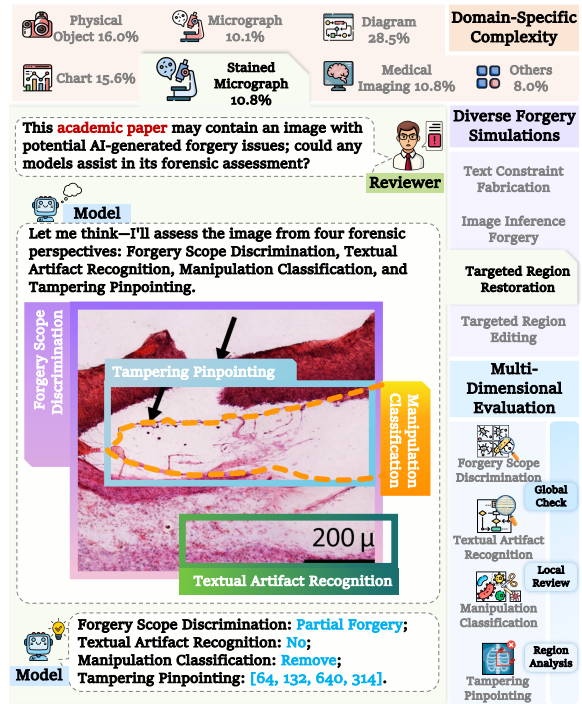


Figure 1: **AEGIS** investigates whether current models can effectively audit *AI-generated images* in academic papers by performing holistic forensic analysis across four complementary tasks.

while, leveraging visual understanding and reasoning, multimodal large language models (MLLMs) have been increasingly applied to image forgery analysis, either directly (Wen et al., 2025) or in conjunction with expert models (Xu et al., 2025).

Despite recent advances, the applicability of current models to expert-level academic review in high-stakes, domain-specific settings remains underexplored. From an expert artificial general intelligence (AGI) perspective (Figure 1), academic image forensics requires models to progress from coarse authenticity judgment to precise region attribution and spatial grounding, as the academic domain presents three specific challenges: (1) more complex visual distributions, (2) finer-grained ma-

Benchmark	Domain	Tasks	# GMs	Forgery Strategies					Evaluation Tasks				# Samples
				TCF	IIF	TRR	TRE	Total	FSD	TAR	MC	TP	
GenImage (NeurIPS’23)		D	6	✓	✓	✗	✗	2	✓	✗	✗	✗	2.7M
Semi-Truths (NeurIPS’24)		D	8	✗	✓	✓	✗	2	✓	✗	✗	✗	1.5M
AIGIBench (NeurIPS’25)		D	16	✓	✗	✗	✗	1	✓	✗	✗	✗	429k
ForensicHub (NeurIPS’25)		D+L	10	✓	✗	✗	✓	2	✓	✗	✗	✓	4M
GIM (AAAI’25)		D+L	3	✗	✗	✗	✓	1	✓	✗	✗	✓	2M
GRE (ACM MM’24)		D+L	3	✗	✗	✗	✓	1	✓	✗	✓	✓	2.3k
DFBench (ACM MM’25)		D	12	✓	✓	✓	✓	4	✓	✗	✗	✗	540k
AIGuard (ACL Findings’25)		D+E	1	✗	✓	✗	✗	1	✓	✗	✗	✗	253k
AEGIS (ours)	Academic	D+L	25	✓	✓	✓	✓	4	✓	✓	✓	✓	20k

Table 1: **Comparison of AEGIS and other benchmarks.** : Nature; : Human; : General; : Daily Life; : News; : E-commerce. **D**: Detection; **L**: Localization; **E**: Explanation. **GMs**: Generative Models (excluding pre-2022 models, which are primarily GAN-based). **TCF**: Text Constraint Fabrication; **IIF**: Image Inference Forgery; **TRR**: Targeted Region Restoration; **TRE**: Targeted Region Editing. **FSD**: Forgery Scope Discrimination; **TAR**: Textual Artifact Recognition; **MC**: Manipulation Classification; **TP**: Tampering Pinpointing.

manipulation patterns, and (3) knowledge-intensive reasoning requirements. This raises the question of *whether current models can generalize their capabilities to support the multi-stage forensic analysis required for academic review*. As shown in Table 1, existing benchmarks fall short of supporting expert-level forensic analysis due to three key limitations:

- **Generic-Scene Bias.** *Academic images are characterized by fine-grained, structured, and semantically dense visual content that poses challenges for forensic interpretation.* In contrast, benchmarks like Semi-Truths (Pal et al., 2024) and ForensicHub (Du et al., 2025) focus on generic imagery such as faces, natural scenes, or social media content with coarser structures, limiting their suitability for academic forensic reasoning.
- **Restricted Forgery Strategies.** *In the context of academic imagery, forgeries exhibit diverse forms with high degrees of subtlety.* However, benchmarks such as AIGuard (Zhang et al., 2025) and GRE (Sun et al., 2024) primarily focus on global forgery or single-type editing forgery, failing to simulate the complex and misleading manipulation behaviors.
- **Forensic-Agnostic Evaluation Protocols.** *Expert academic review requires evidence-grounded forensic assessment beyond single authenticity prediction.* However, benchmarks such as AIGIBench (Li et al., 2025d) and DFBench (Wang et al., 2025) focus on detection accuracy, overlooking key forensic dimensions such as manipulation scope and localization.
- **Domain-Specific Complexity.** **AEGIS** captures the complexity of academic imagery through systematic coverage of **seven major categories and 39 fine-grained subtypes**. Our taxonomy spans diverse scientific visual forms, including microscopic particles, cellular structures, and heterogeneous imaging modalities (e.g., microscopy, medical, and physical imaging), enabling rigorous forensic evaluation in scholarly settings.
- **Diverse Forgery Simulations.** **AEGIS** covers four representative forgery types commonly found in academic publications, including (1) Text Constraint Fabrication, (2) Image Inference Forgery, (3) Targeted Region Restoration, and (4) Targeted Region Editing. These forgeries are instantiated using **25 state-of-the-art generative models**, enabling realistic and adversarial simulation of academic image manipulation.
- **Multi-Dimensional Forensic Evaluation.** **AEGIS** introduces a multi-level and multi-perspective evaluation suite, including (1) Forgery Scope Discrimination, (2) Textual

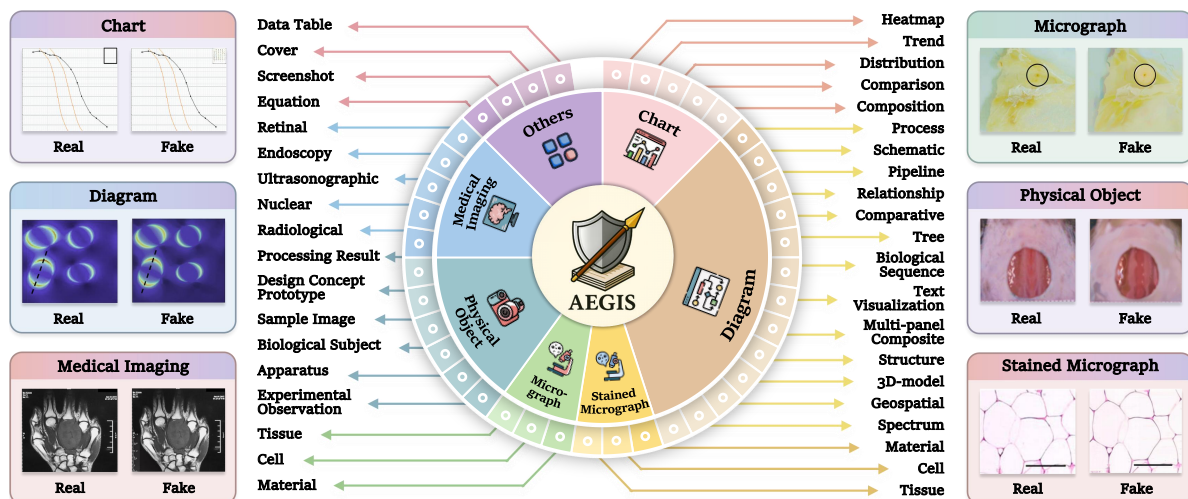


Figure 2: **Hierarchical taxonomy of AEGIS.** We organize academic images into seven categories and 39 fine-grained subtypes based on their structural and semantic characteristics. Real (left) and fake (right) examples are shown for comparison, where “fake” refers to AI-generated forgeries.

Artifact Recognition, (3) Manipulation Classification, and (4) Tampering Pinpointing. These tasks comprehensively assess models’ capabilities in judgment accuracy, localization precision, and explainability depth, closely aligning with real-world academic review needs.

We evaluate 14 proprietary MLLMs, 11 open-source MLLMs, one representative unified multimodal understanding and generation model, as well as six vision-only expert models and three hybrid MLLM-assisted expert models. The experimental results reveal three key findings:

- **Domain-specific complexity challenges detection and localization.** Both MLLMs and expert models struggle on academic images, with GPT-5.1 achieving 48.80% overall performance and expert baselines exhibiting an approximately 30% Real–Forgery F1 gap. This pattern indicates a structural bias toward geometric regularities, failing to generalize to the high-variance, texture-rich nature of academic imagery and leading to severely limited localization accuracy (IoU 30.09%, with only two MLLMs exceeding 50% region-level accuracy).
- **Generative diversity highlights forensic vulnerabilities.** Generative diversity exposes an adversarial asymmetry between generation and forensic analysis. Averaged across models, 11 generative models reduce accuracy below 50%, with four falling below 30%, indicating that existing static forensic cues struggle to generalize against the rapid evolution of generative models.

- **Comprehensive forensic reasoning requires complementarity.** MLLMs excel at forensic reasoning, achieving 84.74% accuracy in textual artifact recognition and 60.07% in manipulation classification, capabilities beyond expert models. In contrast, expert detectors achieve over 70% accuracy in binary authenticity detection but are less robust to post-processing perturbations, indicating that reliable academic image scrutiny requires quantitatively complementary forensic evidence rather than any single paradigm.

2 AEGIS Benchmark

2.1 Overview of AEGIS

We introduce **AEGIS**, a comprehensive benchmark for systematic evaluation of model performance in academic image forensics. As illustrated in Figure 2, **AEGIS** establishes a hierarchical taxonomy with seven main categories (i.e., Chart, Medical Imaging, Physical Object, Micrograph, Stained Micrograph, Diagram, and Others), comprising 39 annotated subtypes. **AEGIS** is constructed based on forgery patterns inspired by genuine retraction cases and academic review observations. By synthesizing evidence from public platforms (e.g., Retraction Watch¹ and PubPeer²) and expert analysis, we identify four recurring forgery strategies. Qualitative expert evaluation, detailed in Appendix A.2, further confirms the consistency between our synthetic data and real-world misconduct cases.

¹<https://retractionwatch.com>

²<https://pubpeer.com>

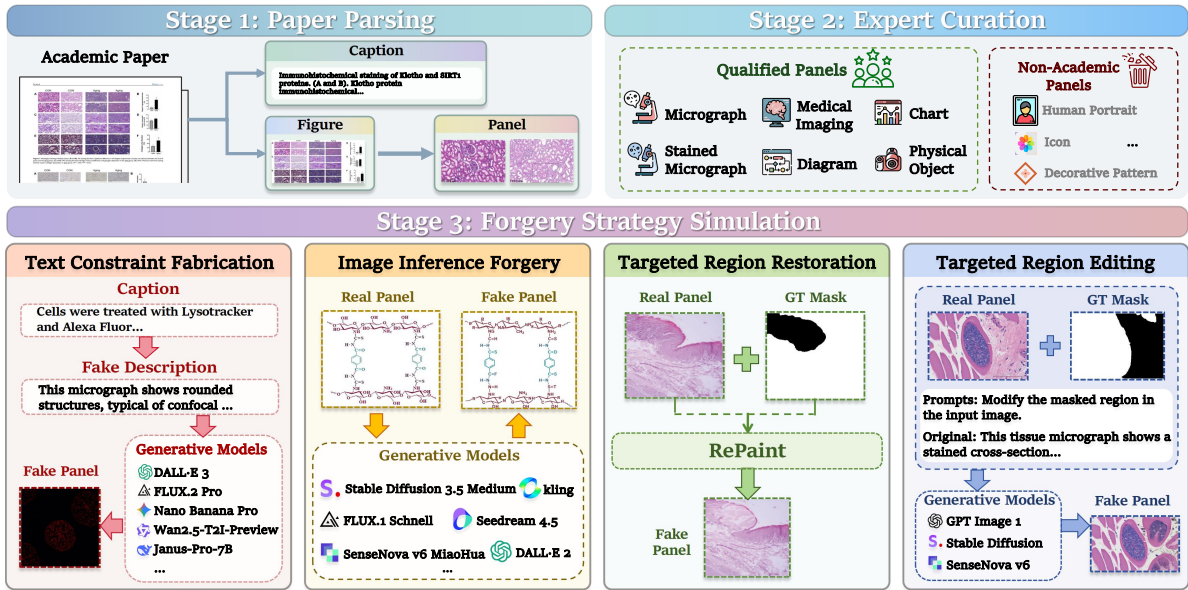


Figure 3: **Construction pipeline of AEGIS. Stage 1: Paper Parsing** extracts figures, captions, and panels from papers. **Stage 2: Expert Curation** retains qualified academic panels while excluding non-academic ones. **Stage 3: Forgery Strategy Simulation** synthesizes AI-generated academic image forgeries via four representative strategies.

2.2 Dataset Construction

Data Curation. Papers were parsed into figures and panels, then manually verified through expert curation to ensure data quality.

- **Paper Parsing.** AEGIS collected more than 4,000 high-quality academic papers from the open-access PubMed Central (PMC) repository³ and performed document-level parsing to extract figures, captions, and panels, where a panel denotes the smallest indivisible visual unit in an academic image. Details for paper selection criteria are in Appendix A.4.
- **Expert Curation and Validation.** We conducted a two-stage quality control process involving manual screening and expert review. Panels with insufficient resolution, parsing errors, or lacking academic semantics were removed, while panels consistent with the predefined hierarchical taxonomy were retained and annotated. Non-academic or non-forgable visual elements (e.g., portraits, icons, and decorative graphics) were excluded, resulting in over 8,000 high-quality annotated panels for dataset construction.

Generative Models. To ensure representative quality and diversity for forensic evaluation, AEGIS incorporates 25 advanced generative models (DeepMind, 2025c; Lugmayr et al., 2022) from

³<https://pmc.ncbi.nlm.nih.gov>

three major architectures: (1) diffusion-based (e.g., Flux, Midjourney V6/V7, DALL-E 2/3, and kling); (2) hybrid (e.g., Wan2.1-T2I-Turbo and GPT Image 1); and (3) unified multimodal understanding and generation (e.g., Janus-Pro-7B (Chen et al., 2025b)). More details are in Appendix A.5.1.

Forgery Strategy Simulations. As illustrated in Figure 3, to simulate realistic forgery behaviors in academic scenarios, AEGIS defines four representative forgery strategies, inspired by common manipulation practices in academic imagery and aligned with the capabilities of contemporary generative models. More details are in Appendix A.5.2.

- **Text Constraint Fabrication (TCF).** TCF models the fabrication of entire academic images from scratch under textual constraints. We reconstructed captions of authentic academic images into semantically equivalent prompts, which were then used to guide text-to-image generative models, resulting in 3,121 Entire Forgeries.
- **Image Inference Forgery (IIF).** IIF generates realistic and detailed globally forged images based on the original input, simulating reference-guided generation. By using authentic images as references, we instructed generative models to synthesize globally consistent forgeries, resulting in 2,274 Entire Forgeries.
- **Targeted Region Restoration (TRR).** TRR simulates the forgery risks associated with recon-

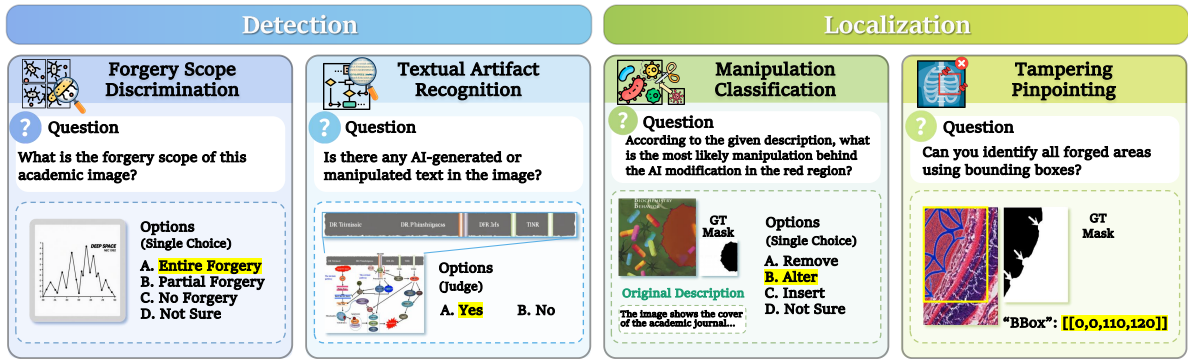


Figure 4: **Evaluation design of AEGIS.** Four evaluation question types are designed to support staged forensic analysis, from global authenticity assessment to fine-grained region-level localization, including **Forgery Scope Discrimination**, **Textual Artifact Recognition**, **Manipulation Classification**, and **Tampering Pinpointing**.

structing missing local content. By applying masks to specific regions of authentic images, we instructed generative models to restore these areas, resulting in 1,650 Partial Forgeries.

- **Targeted Region Editing (TRE).** *TRE emulates common “lightweight tampering” operations.* By providing masks or textual instructions, we performed localized edits on authentic images, including alterations, insertions, or removals, to produce 1,165 Partial Forgeries.

Data Quality Assurance. To mitigate hallucinations from generative models, all AI-generated images underwent a dual-review protocol, evaluating both local manipulation plausibility and global visual fidelity. This process led to the discard of 29% of AI-generated images, resulting in a final set of 8,210 high-quality samples. Verification, performed by five trained annotators, took approximately 200 hours. Detailed assessment is shown in Appendix A.5.3.

2.3 Design of Evaluation Questions

AEGIS designs a set of progressively structured tasks to evaluate forensic capability from global authenticity judgment to fine-grained spatial localization, as shown in Figure 4.

- **Forgery Scope Discrimination (FSD).** This task evaluates coarse-grained authenticity judgment, requiring models to determine whether an input image exhibits no forensic artifacts, contains regionally constrained AI-generated traces, or represents a comprehensively synthesized fabrication. To enable abstention when confidence levels are insufficient, the prompt explicitly incorporates “Not Sure” as an option.

- **Textual Artifact Recognition (TAR).** This task assesses fine-grained visual cue perception by requiring models to perform binary classification on whether textual regions in an input image exhibit traces of AI synthesis, relying on text region recognition and semantic consistency analysis.

- **Manipulation Classification (MC).** This task evaluates structure-aware reasoning, requiring models to infer the type of manipulation applied to an image. Given a red-highlighted manipulated region and the original image caption, the model analyzes the structural and contextual role of the edited area to classify the manipulation as insertion, removal, or alteration.

- **Tampering Pinpointing (TP).** This task focuses on fine-grained spatial grounding, evaluating the localization accuracy of partially edited regions in academic images. We adopt an adaptive granularity protocol: region-level bounding boxes for MLLMs and pixel-level masks for expert models, emphasizing spatial precision beyond global visual cues.

3 Experiments

3.1 Experimental Setup

Benchmarked Models. We evaluate 14 proprietary MLLMs (OpenAI, 2025a,c; DeepMind, 2025a), 11 open-source MLLMs (Li et al., 2025b; Team et al., 2025; AI@Meta, 2025), one unified multimodal understanding and generation model (Chen et al., 2025b), and nine expert models, including six vision-only and three hybrid models, of which only SIDA (Huang et al., 2025) supports multi-granular forgery discrimination.

Model	Detection				Localization				NFI		
	FSD		TAR		MC		TP				
	ACC	M-F1	ACC	M-F1	ACC	M-F1	Region-Level	Pixel-Level			
							CLA	OLR ↓		IoU	F1
Human											
Human	44.20	37.83	76.14	75.02	68.01	51.11	–	–	–	–	
Proprietary MLLMs											
GPT-5.1	50.99	46.16	76.43	76.38	60.07	48.34	46.87	12.25	–	–	48.80
GPT-4.1	66.34	59.98	83.57	82.74	44.55	40.05	25.93	0.14	–	–	43.31
OpenAI o4-mini-high	54.04	49.35	80.71	80.02	36.80	35.54	33.29	5.44	–	–	42.77
Claude Sonnet 4.5	25.36	19.62	58.76	58.70	44.80	39.70	27.41	24.32	–	–	26.83
Gemini 3 Pro Preview	64.37	53.44	84.74	84.67	48.54	45.42	39.14	20.09	–	–	45.79
Gemini 2.5 Flash	60.29	53.78	81.96	80.97	44.80	37.28	47.42	14.07	–	–	47.02
Doubao-Seed-1.6	41.31	35.38	77.23	76.90	50.05	37.08	36.38	0.28	–	–	41.73
Doubao-Seed-1.6-thinking	44.57	37.13	77.92	77.57	58.47	41.85	46.00	0.64	–	–	46.66
Doubao-Seed-1.6-flash	29.26	23.77	60.67	60.39	35.89	30.96	52.36	1.78	–	–	36.03
Doubao-1.5-thinking-vision-pro	39.78	35.31	75.55	75.44	36.65	32.33	47.53	1.71	–	–	42.39
Doubao-1.5-vision-pro	39.87	32.52	72.35	72.32	29.17	28.96	26.96	2.02	–	–	34.07
Qwen3-VL-Plus	38.77	35.72	79.25	78.98	59.28	36.89	5.76	0.53	–	–	16.53
Qwen-VL-Max	43.24	37.30	73.76	73.74	52.06	36.76	40.30	1.32	–	–	43.01
Qwen-VL-Plus	38.20	31.81	51.09	48.23	42.99	32.23	29.37	1.32	–	–	33.89
Open-Source MLLMs											
Gemma 3 27B	44.86	42.36	66.88	58.39	42.49	34.76	48.70	28.56	–	–	37.55
Qwen2.5-VL-72B	37.63	33.16	70.81	70.76	51.76	35.96	32.29	0.53	–	–	38.71
Llama 4 Maverick	39.66	37.00	68.37	68.32	54.62	40.18	14.56	2.88	–	–	29.15
Ministral 3 14B	27.19	25.44	62.41	62.25	53.62	33.89	9.73	12.86	–	–	19.90
LLaVA-NeXT-7B	10.05	12.23	53.82	53.36	59.04	27.01	10.37	6.47	–	–	16.53
LLaVA-NeXT-34B	18.57	17.01	41.71	35.16	32.65	27.19	51.01	0.78	–	–	27.74
LLaVA-Interleave-7B	33.53	20.61	64.68	54.01	26.00	26.67	5.65	8.28	–	–	13.61
LLaVA-Interleave-7B-DPO	35.56	26.09	59.01	53.98	25.15	24.99	25.22	7.50	–	–	28.18
LLaVA-OneVision-7B	39.20	38.76	50.46	50.45	37.80	31.34	3.94	9.98	–	–	11.45
LLaVA-OneVision-72B	30.37	28.57	56.28	55.92	55.31	39.42	33.89	17.05	–	–	33.81
LLaVA-OneVision-7B-Chat	40.54	39.94	50.44	50.26	39.47	32.45	4.14	6.45	–	–	12.19
Unified Multimodal Understanding and Generation Model											
Janus-Pro-7B	24.97	27.18	60.25	42.59	24.18	24.10	5.31	0.25	–	–	13.77
Vision-Only Expert Models											
DRCT (Chen et al., 2024a) *	55.05	39.55	–	–	–	–	–	–	–	–	–
DIRE (Wang et al., 2023b) *	58.16	50.41	–	–	–	–	–	–	–	–	–
AIDE (Yan et al., 2025) *	79.54	44.74	–	–	–	–	–	–	–	–	–
AlignedForensics (Sundara Rajan et al., 2025) *	47.57	42.48	–	–	–	–	–	–	–	–	–
Denoising Trajectory (Liang et al., 2025) *	73.68	65.37	–	–	–	–	–	–	–	–	–
DDA (Chen et al., 2025a) *	74.19	67.07	–	–	–	–	–	–	–	–	–
Hybrid MLLM-Assisted Expert Models											
FakeShield (Xu et al., 2025) *	59.72	62.08	–	–	–	–	–	–	30.09	69.11	–
SIDA (Huang et al., 2025)	52.91	64.38	–	–	–	–	–	–	3.66	54.28	–
FakeVLM (Wen et al., 2025) *	58.73	55.39	–	–	–	–	–	–	–	–	–

Table 2: **Performance on AEGIS across four dimensions.** * denotes binary authenticity classifiers. **FSD**: Forgery Scope Discrimination; **TAR**: Textual Artifact Recognition; **MC**: Manipulation Classification; **TP**: Tampering Pinpointing. **CLA**: Correct Localization Accuracy; **OLR**: Over Localization Rate; **NFI**: Normalized Forensic Index.

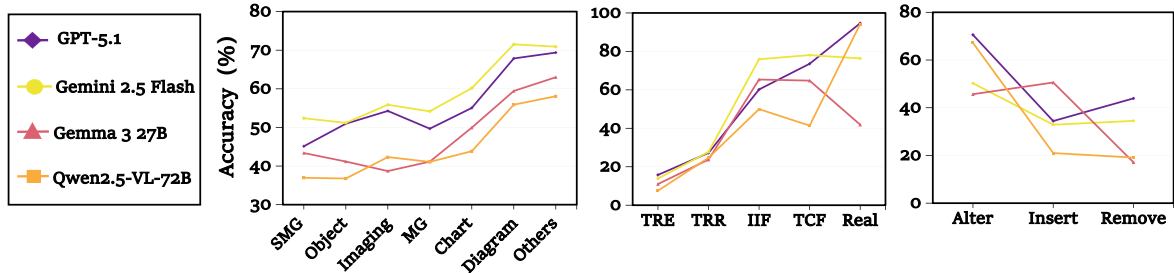
Evaluation Protocol. We uniformly adopt high-resolution PNG images for all evaluations to avoid compression artifacts introduced by lossy formats such as JPEG, which may interfere with forensic cues and bias model judgments.

Evaluation Metrics. We adopt task-specific metrics to systematically evaluate model performance.

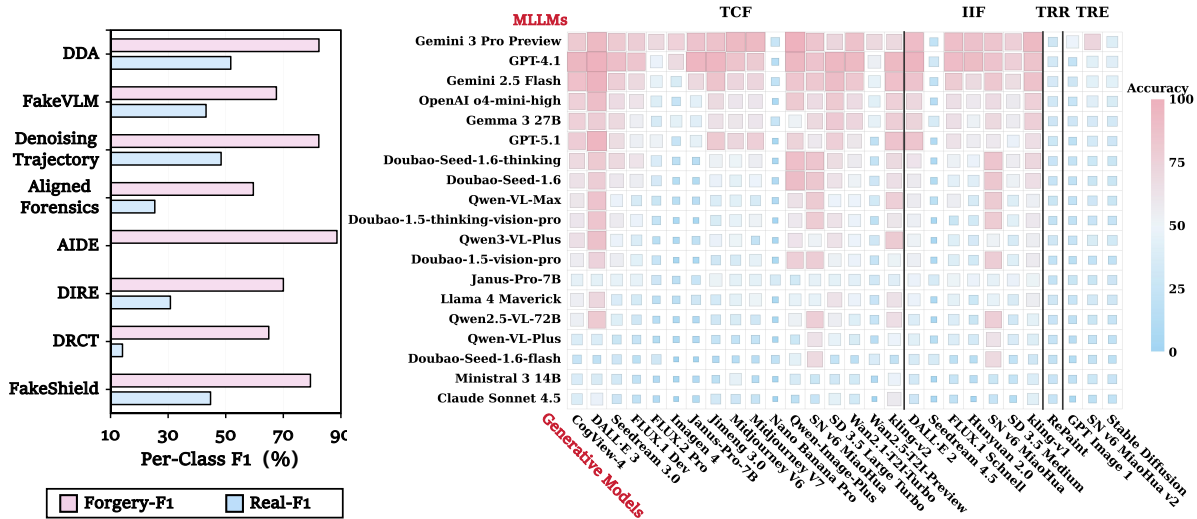
- **Tasks FSD, TAR, and MC.** We employ Accuracy (ACC) and Macro-F1 Score (M-F1) to evaluate performance across tasks.
- **Task TP.** We adopt an adaptive evaluation pro-

cedure: pixel-level metrics, including Intersection over Union (IoU) and F1 Score (F1), evaluate expert-generated masks, while region-level metrics assess MLLM-predicted bounding boxes, where Correct Localization Accuracy (CLA) measures coverage and Over Localization Rate (OLR) penalizes over-localization.

We further propose a new composite metric, the **Normalized Forensic Index (NFI)**, to reflect balanced forensic capability across tasks rather than peak performance on any single dimension. Further details are provided in Appendix B.3.



(a) Impact of Image Categories, Forgery Strategies, and Manipulation Types



(b) Comparison of Per-Class F1 for Expert Models

(c) Adversarial Dynamics Between Generation and Detection

Figure 5: Fine-grained experimental analysis of AEGIS. TCF: Text Constraint Fabrication; IIF: Image Inference Forgery; TRR: Targeted Region Restoration; TRE: Targeted Region Editing. SMG: Stained Micrograph; MG: Micrograph. Seedream: Doubao-Seedream; SD: Stable Diffusion; SN: SenseNova.

3.2 Experimental Analysis

3.2.1 Domain-Specific Complexity Challenges Current Models

We first investigate how the intrinsic complexity of academic imagery impacts forensic capability. Our analysis suggests that the high variance between structured and dense scientific visual forms poses a substantial challenge to the generalization capability of current architectures.

- **A holistic capability gap prevents any model from meeting the multifaceted forensic demands.** As detailed in Table 2, models excel only at isolated tasks, a limitation quantified by the Normalized Forensic Index (NFI) where no evaluated MLLM exceeds 60%, indicating limited balanced forensic capability.
- **Tampering Pinpointing remains a major bottleneck.** Across all expert models, the highest pixel-level IoU reaches only 30.09%, while among MLLMs, only two models exceed 50% CLA.

- **Dense semantics from visual density bias impede generalization.** Performance correlates heavily with visual density. Figure 5(a) shows that robust generalization on structured images (e.g., Chart and Diagram) contrasts sharply with marked degradation on dense categories (e.g., Stained Micrograph and Medical Imaging). This suggests a structural bias: current models tend to over-rely on explicit geometric regularities, while struggling to capture the complex texture details.
- **An authenticity modeling deficit causes expert models to exhibit a pronounced forgery bias.** The challenge of modeling diverse scientific distributions predisposes expert models toward biased prediction strategies. Figure 5(b) shows most expert baselines achieve lower Real-F1 than Forgery-F1. This imbalance suggests that expert models, unable to encompass the manifold of authentic academic imagery, default to flagging unseen patterns as forged.

3.2.2 Generative Diversity Highlights Forensic Vulnerabilities

Adversarial analysis reveals a capability gap: generative evolution outpaces forensic capability, particularly in fine-grained manipulations.

- **The granularity gap makes localized forgeries harder than global ones.** Models perform reliably on globally synthesized images but show notable degradation on localized manipulations. As shown in Figure 5(a), performance on targeted restoration and editing consistently underperforms that on global forgeries, indicating the increased difficulty of localized forgery analysis.
- **Adversarial asymmetry leads generation to outpace forensic capabilities.** Figure 5(c) reveals a clear adversarial asymmetry. Averaged across models, four generative models reduce accuracy below 30% (e.g., Nano Banana Pro), while robust defense remains limited to a few models (e.g., Gemini 3 Pro Preview and GPT-4.1). This result shows that generative advances degrade performance more broadly than forensic capabilities can adapt, exposing a persistent gap in academic image forensics.

3.2.3 Towards Expert AGI: Synergizing Detection and Reasoning

Our findings suggest that no single paradigm is sufficient; instead, the future of academic forensics likely lies in synergizing the sensitivity of experts with the reasoning of MLLMs.

- **Foreground-background entanglement complicates recognition.** Figure 5(a) shows that *insertion* and *removal* are consistently harder to identify than *alteration*. Even the strongest MLLM, GPT-5.1, reaches 60.07% accuracy on Manipulation Classification, indicating reasonable reasoning capability but remaining challenges in separating foreground anomalies from coherent backgrounds.
- **Complementary synergies require orchestrating sensitive sensors with cognitive agents.** The results reveal a functional orthogonality between expert models and MLLMs: experts act as precise but fragile visual “sensors” for authenticity discrimination, while MLLMs provide stronger semantic reasoning (e.g., Textual Artifact Recognition). This contrast is reinforced by their robustness profiles, where vision-centric experts

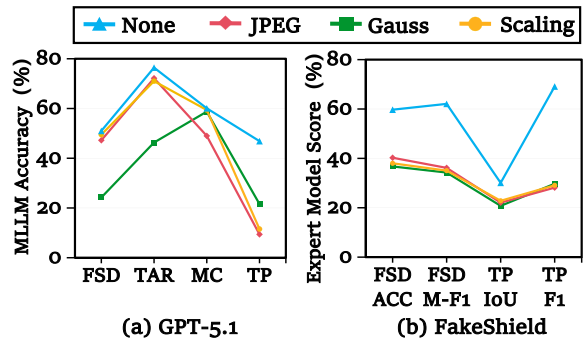


Figure 6: **Impact of post-processing perturbations on AEGIS.** FSD: Forgery Scope Discrimination; TAR: Textual Artifact Recognition; MC: Manipulation Classification; TP: Tampering Pinpointing.

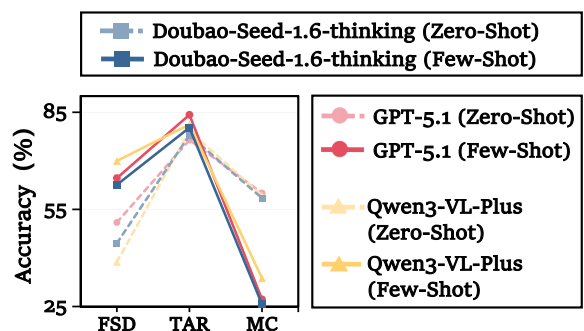


Figure 7: **Impact of Few-Shot prompting on AEGIS.**

degrade sharply under post-processing perturbations (i.e., Gaussian blurring, JPEG compression and image scaling; Figure 6), whereas MLLMs remain comparatively stable, indicating a weaker reliance on low-level visual fidelity. These findings motivate adaptive expert-MLLM orchestration for future **Expert AGI**.

3.3 Prompting Strategies

We evaluate two prompting strategies on AEGIS, namely **Chain-of-Thought** (CoT) (Wei et al., 2022) and **Few-Shot** prompting (Brown et al., 2020). More details are provided in Appendix C.1.

- Few-Shot prompting introduces a clear trade-off (Figure 7): while it enhances pattern-level recognition in detection tasks (e.g., Forgery Scope Discrimination and Textual Artifact Recognition), it consistently degrades fine-grained reasoning required for localized spatial comparison (e.g., Manipulation Classification). This degradation suggests that few-shot examples bias the model toward shortcut pattern matching, which interferes with the multi-step reasoning required to distinguish subtle manipulation types.

- CoT prompting yields consistent gains across most tasks, while substantially degrading manipulation classification. CoT improves detection-centric tasks, indicating that explicit reasoning helps surface more reliable authenticity cues and localization signals. In contrast, Manipulation Classification suffers a significant drop, suggesting limited capability to perform hypothetical judgments.

3.4 Error Analysis

We analyze representative error cases from GPT-5.1 by randomly sampling 100 instances for each task. Full examples and further interpretation are provided in Appendix D.

- **Forgery Scope Discrimination (48%).** Caused by the model’s tendency to overgeneralize localized edits as entire forgeries, or to overlook subtle tampering and misclassify partial forgeries.
- **Textual Artifact Recognition (22%).** Failure to detect fine-grained textual anomalies introduced by text-based modifications in partial forgery.
- **Manipulation Classification (36%).** Stemming from difficulty in distinguishing foreground manipulations from complex background regions.
- **Tampering Pinpointing (56%).** Stemming from a tendency to predict fewer but overly large bboxes, resulting in imprecise localization.

4 Related Work

Existing Benchmarks and Datasets. Prior efforts fall short of supporting expert-level forensic analysis for academic images, due to three structural limitations: (1) *Generic-scene bias*, where datasets such as Semi-Truths (Pal et al., 2024) and GenImage (Zhu et al., 2023) are dominated by faces or natural scenes rather than structured academic images; (2) *Restricted forgery strategies*, as benchmarks including GRE (Sun et al., 2024) and AIGuard (Zhang et al., 2025) largely focus on global or single-type edits and fail to model subtle academic manipulations; and (3) *Forensic-agnostic evaluation protocols*, where datasets such as AIGIBench (Li et al., 2025d) and DFbench (Wang et al., 2025) emphasize detection accuracy while overlooking manipulation scope and localization.

Forensic Models: Expert Models, MLLMs, and Hybrid Approaches. In parallel, the literature has explored a broad spectrum of forensic models for generative forgery analysis, spanning vision-only expert models, MLLMs, and hybrid systems. Vision-only expert models cover multiple detection paradigms, including diffusion-process-based methods such as DIRE (Wang et al., 2023b), DRCT (Chen et al., 2024a), and denoising-trajectory-based methods (Liang et al., 2025), as well as alignment-based approaches including AlignedForensics (Sundara Rajan et al., 2025) and DDA (Chen et al., 2025a). Meanwhile, MLLMs have exhibited strong cross-modal understanding and reasoning capabilities, exemplified by OpenAI o4-mini-high (OpenAI, 2025c), Gemini 2.5 Flash (DeepMind, 2025a), and the Qwen-VL series (Bai et al., 2025b), motivating their evaluation for generative image forensics in complex scenarios. Hybrid approaches integrate MLLMs with expert models either directly, as in FakeVLM (Wen et al., 2025), or in an assisted manner, as exemplified by SIDA (Huang et al., 2025) and FakeShield (Xu et al., 2025).

5 Conclusion

We introduce **AEGIS**, a holistic benchmark comprising over 20k forensic questions, designed to evaluate vision-only models, MLLMs, and hybrid systems for academic image forensics. The primary contribution of **AEGIS** lies in revealing systematic forensic capability gaps that were not exposed by prior benchmarks. These gaps stem from the nature of academic images, which involve fine-grained details, knowledge-intensive content, and structurally complex layouts.

Similar challenges are **not unique to academic imagery**. Domains such as legal evidence analysis (e.g., forged visual exhibits), financial document forensics, and other structured visual materials also require **fine-grained inspection, domain knowledge, and precise region-level reasoning**. These shared characteristics suggest that the capability gaps exposed by **AEGIS** are rooted in common challenges rather than being confined to academic images alone. We hope **AEGIS** provides a robust foundation for advancing trustworthy, interpretable, and generalizable multimodal forensic systems.

Limitations

While **AEGIS** provides a comprehensive benchmark for evaluating AI-generated academic images, several limitations remain:

(1) The availability of *real-world retracted papers involving AI-generated images* is limited. Only a small number of confirmed real cases are publicly accessible, which constrains the scale at which authentic misconduct examples can be incorporated into our benchmark. We expect to include more real-world cases as attention to generative image misuse in academic publishing grows.

(2) Baseline coverage is constrained by the availability of open-source implementations. Although we evaluate a broad range of state-of-the-art models, some recent approaches (e.g., AIGI-Holmes) lack released training weights. Future updates will incorporate additional baselines as more expert models become publicly accessible.

Ethical Considerations

AEGIS complies with ACL ethics guidelines. This study involved no human subjects or animal experimentation. All data were collected from open-source repositories in accordance with relevant usage licenses, ensuring that privacy is preserved and no personally identifiable information is included. **AEGIS** is released under the Creative Commons Attribution 4.0 International License (CC BY 4.0), and the associated codebase is distributed under the Apache License 2.0, supporting both commercial and open-source applications. We have made consistent efforts to minimize bias and ensure transparency throughout the dataset construction and evaluation process.

Acknowledgments

This work is supported by the National Natural Science Foundation of China (Grant Nos. 62473271, 62176026), the Fundamental Research Funds for the Beijing University of Posts and Telecommunications (Grant No. 2025AI4S03), and the BUPT Innovation and Entrepreneurship Support Program (Grant No. 2025-YC-A042). This work is also supported by the Engineering Research Center of Information Networks, Ministry of Education, China. We would also like to thank the anonymous reviewers and area chairs for constructive discussions and feedback.

References

- AI@Meta. 2025. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>. Accessed: 2025-04-05.
- Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, and 45 others. 2025a. **Qwen3-vl technical report**. *Preprint*, arXiv:2511.21631.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025b. **Qwen2.5-vl technical report**. *Preprint*, arXiv:2502.13923.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. **Language models are few-shot learners**. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- ByteDance. 2025a. Doubao-1.5-thinking-pro model card. <https://console.volcengine.com/ark/region:ark+cn-beijing/model/detail?Id=doubao-1-5-thinking-pro>. Accessed: 2025-04-15.
- ByteDance. 2025b. Doubao-1.5-thinking-vision-pro model card. <https://console.volcengine.com/ark/region:ark+cn-beijing/model/detail?Id=doubao-1-5-thinking-vision-pro>. Accessed: 2025-04-28.
- ByteDance. 2025c. doubao-seed-1.6-flash model card. <https://console.volcengine.com/ark/region:ark+cn-beijing/model/detail?Id=doubao-seed-1-6-flash>. Accessed: 2025-06-11.
- ByteDance. 2025d. doubao-seed-1.6-thinking model card. <https://console.volcengine.com/ark/region:ark+cn-beijing/model/detail?Id=doubao-seed-1-6-thinking>. Accessed: 2025-06-11.
- Baoying Chen, Jishen Zeng, Jianquan Yang, and Rui Yang. 2024a. **DRCT: Diffusion reconstruction contrastive training towards universal detection of diffusion generated images**. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 7621–7639. PMLR.

- Jiaxuan Chen, Jieteng Yao, and Li Niu. 2024b. **A single simple patch is all you need for ai-generated image detection**. *Preprint*, arXiv:2402.01123.
- Ruoxin Chen, Junwei Xi, Zhiyuan Yan, Ke-Yue Zhang, Shuang Wu, Jingyi Xie, Xu Chen, Lei Xu, Isabel Guan, Taiping Yao, and Shouhong Ding. 2025a. **Dual data alignment makes AI-generated image detector easier generalizable**. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. 2025b. **Janus-pro: Unified multimodal understanding and generation with data and model scaling**. *Preprint*, arXiv:2501.17811.
- Riccardo Corvi, Davide Cozzolino, Giada Zingarini, Giovanni Poggi, Koki Nagano, and Luisa Verdoliva. 2023. **On the detection of synthetic images generated by diffusion models**. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.
- Google DeepMind. 2025a. Gemini 2.5 flash and pro are now generally available, and we’re introducing 2.5 flash-lite, our most cost-efficient and fastest 2.5 model yet. <https://blog.google/products/gemini/gemini-2-5-model-family-expands/>. Accessed: 2025-06-17.
- Google DeepMind. 2025b. Gemini 3 pro model card. <https://storage.googleapis.com/deepmind-media/Model-Cards/Gemini-3-Pro-Model-Card.pdf>. Accessed: 2026-01-02.
- Google DeepMind. 2025c. Imagen 4 generate preview 06-06. <https://cloud.google.com/vertex-ai/generative-ai/docs/models/imagen/4-0-generate-preview-06-06?hl=zh-cn>. Accessed: 2025-06-06.
- Bo Du, Xuekang Zhu, Xiaochen Ma, Chenfan Qu, Kaiwen Feng, Zhe Yang, Chi-Man Pun, Jian liu, and Ji-Zhe Zhou. 2025. **Forensichub: A unified benchmark & codebase for all-domain fake image detection and localization**. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. **CLIPScore: A reference-free evaluation metric for image captioning**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zhenglin Huang, Jinwei Hu, Xiangtai Li, Yiwei He, Xingyu Zhao, Bei Peng, Baoyuan Wu, Xiaowei Huang, and Guangliang Cheng. 2025. **Sida: Social media image deepfake detection, localization and explanation with large multimodal model**. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 28831–28841.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollar, and Ross Girshick. 2023. **Segment anything**. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4015–4026.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2025a. **LLaVA-onevision: Easy visual task transfer**. *Transactions on Machine Learning Research*.
- Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun MA, and Chunyuan Li. 2025b. **Llava-interleave: Tackling multi-image, video, and 3d in large multimodal models**. In *International Conference on Learning Representations*, volume 2025, pages 81182–81199.
- Yumeng Li, Guang Yang, Hao Liu, Bowen Wang, and Colin Zhang. 2025c. **dots.ocr: Multilingual document layout parsing in a single vision-language model**. *Preprint*, arXiv:2512.02498.
- Ziqiang Li, Jiazhen Yan, Ziwen He, Kai Zeng, Weiwei Jiang, Lizhi Xiong, and Zhangjie Fu. 2025d. **Is artificial intelligence generated image detection a solved problem?** In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Yachao Liang, Min Yu, Gang Li, Jianguo Jiang, Fuqiang Du, Li Jingyuan, Lanchi Xie, Zhen Xu, and Weiqing Huang. 2025. **Denoising trajectory biases for zero-shot AI-generated image detection**. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024. **Llava-next: Improved reasoning, ocr, and world knowledge**. <https://llava-vl.github.io/blog/2024-01-30-llava-next/>. Accessed: 2025-04-05.
- Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. 2022. **Repaint: Inpainting using denoising diffusion probabilistic models**. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11461–11471.
- OpenAI. 2025a. Introducing gpt-4.1 in the api. <https://openai.com/index/gpt-4-1/>. Accessed: 2025-04-14.
- OpenAI. 2025b. Introducing gpt-5.1 in the api. <https://openai.com/index/gpt-5-1/>. Accessed: 2025-11-12.
- OpenAI. 2025c. Openai o3 and o4-mini system card. <https://openai.com/index/o3-o4-mini-system-card/>. Accessed: 2025-04-16.

- Anisha Pal, Julia Kruk, Mansi Phute, Manognya Bhataram, Diyi Yang, Duen Horng Chau, and Judy Hoffman. 2024. [Semi-truths: A large-scale dataset of ai-augmented images for evaluating robustness of ai-generated image detectors](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 118025–118051. Curran Associates, Inc.
- Maximilian Seitzer. 2020. pytorch-fid: FID Score for PyTorch. <https://github.com/mseitzer/pytorch-fid>. Accessed: 2025-04-08.
- Oriane Siméoni, Huy V. Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, Francisco Massa, Daniel Haziza, Luca Wehrstedt, Jianyuan Wang, Timothée Darcet, Théo Moutakanni, Leonel Sentana, Claire Roberts, Andrea Vedaldi, and 7 others. 2025. [DINOv3](#). Preprint, arXiv:2508.10104.
- Zhihao Sun, Haipeng Fang, Juan Cao, Xinying Zhao, and Danding Wang. 2024. [Rethinking image editing detection in the era of generative ai revolution](#). In *Proceedings of the 32nd ACM International Conference on Multimedia*, MM '24, page 3538–3547, New York, NY, USA. Association for Computing Machinery.
- Anirudh Sundara Rajan, Utkarsh Ojha, Jedidiah Schloesser, and Yong Jae Lee. 2025. [Aligned datasets improve detection of latent diffusion-generated images](#). In *International Conference on Learning Representations*, volume 2025, pages 63304–63322.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025. [Gemma 3 technical report](#). Preprint, arXiv:2503.19786.
- Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. 2023a. [Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7464–7475.
- Jiarui Wang, Huiyu Duan, Juntong Wang, Ziheng Jia, Woo Yi Yang, Xiaorong Zhu, Yu Zhao, Jiaying Qian, Yuke Xing, Guangtao Zhai, and Xiongkuo Min. 2025. [Dfbench: Benchmarking deepfake image detection capability of large multimodal models](#). In *Proceedings of the 33rd ACM International Conference on Multimedia*, MM '25, page 12666–12673, New York, NY, USA. Association for Computing Machinery.
- Zhendong Wang, Jianmin Bao, Wengang Zhou, Weilun Wang, Hezhen Hu, Hong Chen, and Houqiang Li. 2023b. [Dire for diffusion-generated image detection](#). In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22445–22455.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.
- Siwei Wen, Junyan Ye, Peilin Feng, Hengrui Kang, Zichen Wen, Yize Chen, Jiang Wu, wenjun wu, Conghui He, and Weijia Li. 2025. [Spot the fake: Large multimodal model-based synthetic image detection with artifact explanation](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Zhipei Xu, Xuanyu Zhang, Runyi Li, Zecheng Tang, Qing Huang, and Jian Zhang. 2025. [Fakeshield: Explainable image forgery detection and localization via multi-modal large language models](#). In *International Conference on Learning Representations*, volume 2025, pages 31186–31216.
- Shilin Yan, Ouxiang Li, Jiayin Cai, Yanbin Hao, Xiaolong Jiang, Yao Hu, and Weidi Xie. 2025. [A sanity check for ai-generated image detection](#). In *International Conference on Learning Representations*, volume 2025, pages 70702–70720.
- Wenhua Zhang, Weicheng Li, Xuanrong Rao, Lixin Zou, Xiangyang Luo, Chubin Zhuang, Yongjie Hong, Zhen Qin, Hengyu Chang, Chenliang Li, and Bo Zheng. 2025. [AIGuard: A benchmark and lightweight detection for E-commerce AIGC risks](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 12437–12450, Vienna, Austria. Association for Computational Linguistics.
- Mingjian Zhu, Hanting Chen, Qiangyu YAN, Xudong Huang, Guanyu Lin, Wei Li, Zhijun Tu, Hailin Hu, Jie Hu, and Yunhe Wang. 2023. [Genimage: A million-scale benchmark for detecting ai-generated image](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 77771–77782. Curran Associates, Inc.

attributed to AI-generated academic images. This scarcity does not indicate that AI-generated image misuse is rare in practice, but rather reflects the inherent difficulty of detecting such forgeries in real-world scientific publishing.

Practical Constraints of Using Real-World Misconduct Cases. Currently available and verifiable misconduct cases involving AI-generated academic images exhibit several limitations:

- **Lack of Structured Annotations.** Most cases do not provide region- or pixel-level ground truth, which prevents reliable computation of localization metrics and undermines evaluation reproducibility.
- **Legal and Ethical Sensitivity.** Some cases involve ongoing investigations or sensitive ethical issues, making them unsuitable for inclusion in a publicly released benchmark.
- **Lack of Editing Provenance.** Real cases rarely disclose detailed editing procedures or generation sources (e.g., generative models), making it impossible to control manipulation granularity or construct fully verified ground truth.
- **Insufficient Sample Size.** The number of currently available cases is too limited to support statistically stable and systematic evaluation.

Under these conditions, incorporating real-world cases into a systematic benchmark would compromise controllability, reproducibility, and consistency in manipulation granularity. Therefore, **AEGIS** adopts a structured modeling framework grounded in forgery patterns identified by academic review experts, enabling controlled, fine-grained, and reproducible evaluation.

A.2 Qualitative Expert Evaluation

To validate the consistency between synthetic samples and real-world misconduct cases, we conducted a blinded expert evaluation on 200 images randomly sampled from the combined pool of real-world misconduct cases and **AEGIS** synthetic samples. All available authentic cases were included, with the remaining images randomly drawn from the synthetic set. Five experts reviewed the images without knowing their origin. Each image was rated on a 5-point Likert scale (1 = very low, 5 = very high) based on the following criteria:

- **Visual Realism** (how natural the image appears)

Criterion	Real-World	Synthetic
Visual Realism \uparrow	3.86 ± 0.52	4.21 ± 0.56
Structural Plausibility \uparrow	3.98 ± 0.47	4.09 ± 0.50
Consistency with Misconduct Patterns \uparrow	4.18 ± 0.60	4.15 ± 0.63
Overall Credibility \uparrow	4.05 ± 0.54	4.19 ± 0.58

Table 3: **Qualitative expert evaluation.** Comparison between *real-world misconduct cases* and *synthetic samples*. Higher scores indicate better performance.

- **Structural Plausibility** (whether the image structure is coherent and scientifically reasonable)
- **Consistency with Known Misconduct Patterns** (whether the forgery resembles patterns observed in real-world retraction cases)
- **Overall Credibility** (as an academic image)

Table 3 shows that **AEGIS** synthetic samples received scores highly comparable to real-world retraction cases across all four criteria. These findings suggest that **the synthetic data reasonably aligns with real-world misconduct cases in terms of visual realism, structural plausibility, and manipulation characteristics.**

A.3 Dataset Statistics

As shown in Figure 10, we summarize the data distribution of **AEGIS**, comprising real images (1,795 images) and four forgery strategies: Text Constraint Fabrication (3,121 images), Image Inference Forgery (2,274 images), Targeted Region Restoration (1,650 images), and Targeted Region Editing (1,165 images).

As shown in Figure 11, the seven academic image categories exhibit a deliberately diverse distribution, reflecting broad coverage of structured, dense, and heterogeneous scientific visual forms.

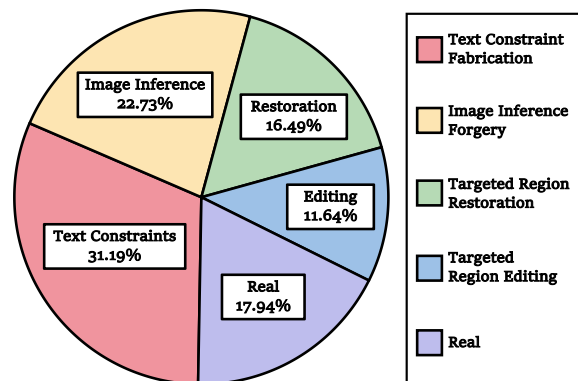


Figure 10: **Distribution of real images and images generated by four forgery strategies in AEGIS.**

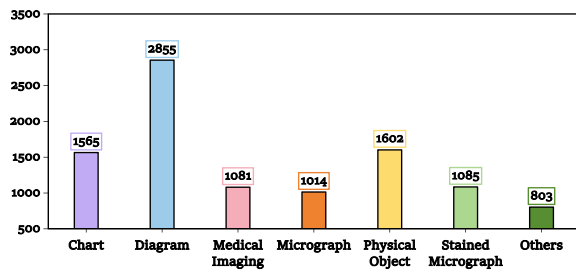


Figure 11: Distribution of seven academic image categories in AEGIS.

A.4 Data Curation

Paper Parsing. AEGIS collected 4,362 high-quality academic papers from the open-access PMC repository and performed document-level parsing to construct the initial visual corpus. The paper selection criteria were as follows:

- **Structural Completeness.** Each paper must contain at least four independent figures to ensure sufficient visual information density. We prioritized papers with multi-panel figures, where a single figure can be reasonably segmented into multiple panels to support subsequent panel-level forgery simulation.
- **Image Parsability and Resolution Requirements.** Figures must maintain adequate clarity in the original publication. Papers were excluded if all images had a shortest side of 100 pixels or less, or exhibited severe scanning blur or other distortions that hindered reliable parsing.
- **Data Cleaning and Deduplication.** We removed duplicated papers and excluded documents lacking structural value, such as those containing only decorative or non-informative images.

As illustrated in Figure 12, a *figure* refers to a complete figure unit paired one-to-one with its corresponding caption, whereas a *panel* denotes the smallest indivisible visual unit within an academic image. The entire extraction and parsing process required approximately 200 hours of human and computational effort.

The extraction and parsing pipeline consisted of the following steps:

- **Step 1: Figure and Caption Extraction.** We employed the Fitz library to extract figures from PDF documents, as it consistently preserved higher visual fidelity than alternative document parsing tools, which was critical for subsequent

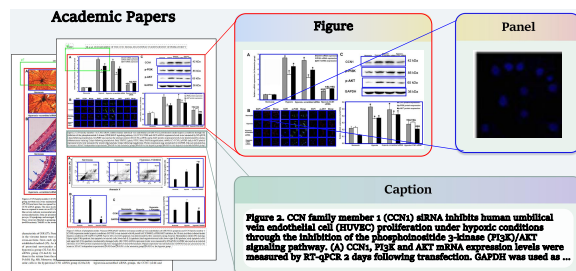


Figure 12: Extraction and Parsing of academic papers. The regions highlighted with red boxes correspond to figures, while those highlighted with blue boxes correspond to panels. Green boxes are used to verify whether the visualization is effective.

panel-level analysis. Captions were parsed using dots.ocr (Li et al., 2025c) and were further used to cross-validate the correctness of the extracted figures. Finally, three human experts with academic peer-review experience manually verified the figure–caption correspondences to ensure extraction accuracy.

- **Step 2: Panel Segmentation from Figures.** We applied YOLOv7 (Wang et al., 2023a) to segment figures into panels. The YOLOv7 model was fully trained on 200,000 manually annotated academic panels and had been validated to meet production-level reliability and commercial deployment standards. Panels with irregular layouts, poor image quality, or lacking experimental visual content were subsequently filtered out through joint inspection by three expert reviewers and five trained graduate annotators.

A.5 Synthetic Data Generation

As discussed in the limitations, real-world misconduct cases are inherently scarce and present multiple practical constraints for benchmark construction. Therefore, we adopted a synthetic data generation paradigm, which enabled scalable, controllable, and well-annotated construction of forensic samples. This subsection describes the synthetic data generation pipeline employed in AEGIS.

A.5.1 Generative Models

As shown in Table 5, we incorporate 25 state-of-the-art generative models spanning diffusion-based models, hybrid architectures, and one unified multimodal understanding and generation model to support diverse academic forgery simulations.

A.5.2 Forgery Simulations

To systematically simulate realistic manipulations in academic imagery, **AEGIS** adopted a strategy-driven forgery simulation framework. Inspired by common misconduct patterns in scholarly publications and aligned with the capabilities of contemporary generative models, we define four representative forgery strategies: *Text Constraint Fabrication*, *Image Inference Forgery*, *Targeted Region Restoration*, and *Targeted Region Editing*. These strategies cover global and localized manipulation scenarios, ranging from semantic-level content fabrication to fine-grained region-specific edits.

Specifically, the first two strategies focus on holistic image synthesis conditioned on textual or contextual constraints, yielding globally generated academic figures. For *Text Constraint Fabrication*, we instructed GPT-4o mini to semantically reconstruct the captions of authentic academic images.

In contrast, the latter two strategies involve localized manipulations applied to selected regions within otherwise authentic images. For these region-level forgeries, we employed the Segment Anything Model (SAM) (Kirillov et al., 2023) to automatically generate precise manipulation masks, ensuring accurate spatial control over the forged regions and enabling consistent pixel-level ground-truth annotations for localization evaluation.

A.5.3 Data Quality Assurance

At the local level, annotators examined whether manipulated regions exhibited coherent textures, well-defined boundaries, and semantic consistency with surrounding content, while avoiding trivial artifacts such as unnatural edges, repetitive patterns, or mismatched visual semantics. At the global level, images were assessed for overall structural coherence, adherence to scientific conventions, and consistency with the corresponding academic context, including figure layout, visual hierarchy, and modality-specific characteristics. The review process involved five human experts with prior academic peer-review experience, together with three Ph.D. students and six trained Master’s students.

In addition to expert review, we conducted an automated quality assessment to contextualize the realism and semantic fidelity of the synthesized images relative to datasets used by baseline forensic models. As shown in Table 4, we compared the quality scores of **AEGIS** against baseline datasets using three metrics: Inception Score (IS), Fréchet Inception Distance (FID), and CLIP Score.

Dataset	Score		
	CLIP \uparrow	FID \downarrow	IS \uparrow
DiffusionForensics	–	27.46	13.47 \pm 0.97
DRCT-2M	75.74	22.97	34.57 \pm 1.57
SID-Set (Full_Synthetic)	–	60.02	38.73 \pm 0.95
SID-Set (Tampered)	–	30.20	32.28 \pm 2.39
New Generator	–	75.58	24.02 \pm 1.35
Robust_LDM	–	99.80	13.76 \pm 1.02
FakeClue	70.95	86.59	13.30 \pm 1.69
EvalGEN	–	69.60	29.67 \pm 0.67
Chameleon	–	53.88	16.87 \pm 0.46
Ours	78.66	32.22	13.41 \pm 1.23

Table 4: Comparison of different datasets.

- **Inception Score (IS)**. IS measures visual quality and content diversity of generated images based on the entropy of class predictions produced by a pretrained Inception network. Higher scores generally indicate images that are both visually coherent and diverse. We note that **IS was originally proposed for evaluating GAN-based image generators, and its sensitivity may be reduced for modern non-GAN architectures** (e.g., diffusion-based or hybrid generative models), whose outputs often depart from the ImageNet-centric class distributions assumed by the Inception model. As a result, lower IS values in our setting should not be interpreted as diminished perceptual quality, but rather as a known limitation of the metric when applied beyond its original design scope.
- **Fréchet Inception Distance (FID)**. FID (Seitzer, 2020) evaluates the distributional similarity between generated images and real academic images by computing the Fréchet distance between their feature embeddings extracted from a pretrained Inception model. Lower FID scores indicate higher visual fidelity and closer alignment with real-image statistics.
- **CLIP Score**. CLIP Score (Hessel et al., 2021) measures cross-modal semantic consistency between generated images and their corresponding textual prompts. It is computed as the cosine similarity between visual and textual embeddings extracted by a pretrained CLIP model, where higher scores reflect stronger semantic alignment.

Together, these complementary metrics provide an objective reference for evaluating perceptual fidelity, distributional realism, and semantic alignment, supporting the validity of the curated dataset alongside expert-driven quality control.

B Experiment Details

B.1 Experiment Environment

For evaluation experiments, most model inferences were conducted via the OpenRouter API⁴. Exceptions include the LLaVA series, which were downloaded from Hugging Face⁵ and executed locally, as well as the Doubao series, which were accessed via the Volcano Engine API⁶.

The system configuration is summarized below:

- **CPU:** Dual-socket Intel Xeon Gold 6148 (2.40 GHz), 20 cores per socket, 80 threads total
- **GPU:** 8 × NVIDIA A40 (48 GB VRAM each)
- **GPU Driver:** 575.57.08
- **CUDA:** 11.8
- **cuDNN:** 8.9.6 (compiled with CUDA 11.8)
- **Operating System:** Ubuntu 22.10

B.2 Benchmarked Models

As summarized in Table 6, we evaluate nine expert models in total, including six vision-only expert models and three hybrid expert models.

As shown in Table 7, we evaluate 25 MLLMs (consisting of 14 proprietary and 11 open-source models) as well as one unified multimodal understanding and generation model. The proprietary models are released by OpenAI (OpenAI, 2025a,b,c), Google (DeepMind, 2025a,b), Alibaba (Bai et al., 2025a), and ByteDance (ByteDance, 2025b,a,c,d). The open-source models are released by Meta (AI@Meta, 2025), Mistral, Alibaba (Bai et al., 2025b), Google (Team et al., 2025), and the LLaVA community (Li et al., 2025a,b; Liu et al., 2024).

B.3 Evaluation Metrics

To comprehensively assess model performance across different dimensions, we design **seven core metrics**.

(1) Accuracy (ACC). ACC measures the proportion of correctly predicted answers among all evaluated instances. All tasks evaluated with ACC are formulated as *single-choice classification problems*, where each image corresponds to exactly one ground-truth label and the model selects a single option.

⁴<https://openrouter.ai>

⁵<https://huggingface.co>

⁶<https://www.volcengine.com>

Let N denote the total number of valid evaluation samples, and let $\mathbb{I}(\cdot)$ be the indicator function. ACC is computed as:

$$\text{ACC} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\hat{y}_i = y_i), \quad (1)$$

where \hat{y}_i and y_i denote the predicted and ground-truth labels of the i -th sample, respectively.

For the Forgery Scope Discrimination (FSD) task, ACC is reported under two settings depending on model capability. Models that support fine-grained scope prediction are evaluated under a multi-class setting (*Real*, *Entire Forgery*, *Partial Forgery*). Vision-only expert models that only support authenticity detection are evaluated under a binary setting (*Real*, *Forgery*), where both global and local forgeries are mapped to the *Forgery* class.

(2) Macro-F1 Score (M-F1). M-F1 is adopted to measure class-balanced performance by averaging F1 scores across classes with equal weight. Similar to ACC, M-F1 is computed over single-choice classification outputs and mitigates the influence of class imbalance.

For each class $c \in \mathcal{C}$, precision P_c , recall R_c , and F1 Score $F1_c$ are defined as:

$$P_c = \frac{\text{TP}_c}{\text{TP}_c + \text{FP}_c}, \quad R_c = \frac{\text{TP}_c}{\text{TP}_c + \text{FN}_c}, \quad (2)$$

$$F1_c = \frac{2 \cdot P_c \cdot R_c}{P_c + R_c}, \quad (3)$$

where TP_c , FP_c , and FN_c denote the numbers of true positives, false positives, and false negatives for class c , respectively.

M-F1 is then computed as:

$$\text{M-F1} = \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} F1_c. \quad (4)$$

For the FSD task, M-F1 is computed separately for the binary and multi-class settings, consistent with the corresponding ACC evaluation protocol. Binary M-F1 is reported for expert models distinguishing *Real* versus *Forgery*, while multi-class M-F1 is reported for models predicting *Real*, *Entire Forgery*, and *Partial Forgery*. Samples with abstention responses (i.e., *Not Sure*) are excluded from M-F1 computation.

(3) Correct Localization Accuracy (CLA).

CLA is a region-level localization metric that evaluates whether all ground-truth tampered regions are sufficiently covered by the predicted regions. CLA is computed at the *sample level* and then aggregated across samples.

For a given sample, let N_{gt} denote the number of ground-truth connected components, and let GT_i denote the i -th ground-truth component. Let Pred represent the union of all predicted regions. A sample is considered *CLA-correct* if and only if every ground-truth component is sufficiently covered by the predicted regions:

$$\text{CLA}_{\text{correct}} = \begin{cases} \text{False}, & \text{if } N_{\text{gt}} = 0, \\ \text{True}, & \text{if } \forall i \in [1, N_{\text{gt}}], \frac{|\text{Pred} \cap \text{GT}_i|}{|\text{GT}_i|} \geq \tau, \\ \text{False}, & \text{otherwise,} \end{cases} \quad (5)$$

where $|\cdot|$ denotes the number of pixels in a region, and τ is a coverage threshold (set to 0.5 in all experiments).

The final CLA score is reported as the percentage of CLA-correct samples over all evaluated samples:

$$\text{CLA} (\%) = \frac{\#\{\text{CLA-correct samples}\}}{\#\{\text{total samples}\}} \times 100\%. \quad (6)$$

(4) Over Localization Rate (OLR). OLR measures whether a model produces excessively many predicted regions relative to the ground-truth, capturing the tendency of over-fragmented or overly sensitive localization.

For a given sample, let N_{gt} denote the number of ground-truth connected components and N_{pred} denote the number of predicted regions. A sample is considered *over-localized* if:

$$\text{OLR}_{\text{over}} = \begin{cases} \text{True}, & \text{if } N_{\text{gt}} = 0 \text{ and } N_{\text{pred}} > 0, \\ \text{True}, & \text{if } N_{\text{gt}} > 0 \text{ and } N_{\text{pred}} > N_{\text{gt}} \times (1 + \alpha), \\ \text{False}, & \text{otherwise,} \end{cases} \quad (7)$$

where α is an over-localization tolerance factor (set to 0.5 by default).

The final OLR score is computed as the proportion of over-localized samples among all evaluated samples:

$$\text{OLR} (\%) = \frac{\#\{\text{over-localized samples}\}}{\#\{\text{total samples}\}} \times 100\%. \quad (8)$$

(5) Intersection over Union (IoU). For pixel-level localization tasks, we adopt IoU to measure the spatial overlap between predicted tampering

masks and ground-truth masks. IoU evaluates localization accuracy at the pixel granularity and is widely used for segmentation-based forensic analysis.

Let Pred and GT denote the predicted binary mask and the ground-truth binary mask for a given sample, respectively. IoU is defined as:

$$\text{IoU} = \frac{|\text{Pred} \cap \text{GT}|}{|\text{Pred} \cup \text{GT}|}, \quad (9)$$

where $|\cdot|$ denotes the number of pixels in the corresponding region.

The final IoU score is reported as the average IoU over all evaluated samples.

(6) Pixel-level F1 Score (F1). To complement IoU, we further report the pixel-level F1 Score, which balances precision and recall of predicted tampering masks. This metric emphasizes accurate boundary delineation and penalizes both over- and under-segmentation.

At the pixel level, true positives (TP), false positives (FP), and false negatives (FN) are defined based on pixel-wise correspondence between the predicted mask and the ground-truth mask. Precision P , recall R , and F1 Score are computed as:

$$P = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad R = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad (10)$$

$$\text{F1} = \frac{2 \cdot P \cdot R}{P + R}. \quad (11)$$

The reported pixel-level F1 Score is averaged across all samples in the evaluation set.

(7) Normalized Forensic Index (NFI). Beyond task-specific metrics, we define a composite score to holistically measure a model's overall forensic capability with an emphasis on balanced performance. NFI is designed to reward models that perform consistently across heterogeneous forensic tasks, while penalizing pathological localization behaviors such as excessive over-localization.

For each model i , we first collect its core forensic scores:

$$\mathcal{S}_i = \{s_{i,1}, s_{i,2}, s_{i,3}, s_{i,4}\}, \quad (12)$$

where \mathcal{S}_i includes the FSD-F1, TAR-F1, MC-F1, and CLA scores of model i . All scores are normalized to $[0, 1]$ for numerical consistency.

We compute the harmonic mean to emphasize balanced performance:

$$\text{HM}_i = \frac{|\mathcal{S}_i|}{\sum_{s \in \mathcal{S}_i} \frac{1}{s + \epsilon}}, \quad (13)$$

where ϵ is a small constant to ensure numerical stability.

To penalize over-localized predictions, we introduce an over-localization factor based on the OLR:

$$P_i = (1 - \text{OLR}_i)^\gamma, \quad (14)$$

where $\text{OLR}_i \in [0, 1]$ and γ controls the penalty strength.

Finally, the NFI is defined as:

$$\text{NFI}_i = 100 \cdot \text{HM}_i \cdot P_i, \quad \gamma = 0.5. \quad (15)$$

A higher NFI indicates stronger and more reliable forensic capability, reflecting both balanced task performance and structurally reasonable localization behavior.

Sensitivity Analysis of NFI. To examine the robustness of the proposed Normalized Forensic Index (NFI), we conducted a sensitivity analysis with respect to the OLR penalty exponent γ , as reported in Table 8.

Table 8 shows that varying γ within a reasonable range ($\gamma \in \{0.3, 0.5, 0.7\}$) leads to only minor changes in absolute NFI values, while the relative ranking of top-performing models remains largely stable. In particular, leading proprietary MLLMs consistently occupy the top ranks, indicating that NFI is not overly sensitive to the exact strength of the OLR-based penalty.

B.4 Evaluation Protocol

We uniformly adopted **high-resolution PNG images** for all evaluations to avoid compression artifacts introduced by lossy formats (e.g., JPEG), which may interfere with forensic cues and bias model judgments. When an input image exceeded a model’s maximum input resolution, a two-stage adaptive compression strategy was applied: (1) the long edge was first scaled to the model-specific maximum input size while preserving the aspect ratio; and (2) the total pixel count was further reduced, if necessary, through iterative 10% down-scaling until the model’s input constraints were satisfied.

All prompts were intentionally kept minimal, containing only task definitions and direct instructions, in order to reduce extraneous noise and better expose the intrinsic forensic capabilities of the evaluated models.

C Multi-Dimensional Evaluation Analysis

C.1 Impact of Prompting Strategies

C.1.1 Impact of Few-Shot Prompting

For Few-Shot prompting, we built an indexed reference library of curated examples. Given an input image, retrieval was restricted to the same category based on dataset annotations. Visual similarity was then computed using DINO3 (Siméoni et al., 2025) embeddings, and the single most similar reference image was selected as the in-context exemplar. The user prompt included the image description, the input image, and an output-format example to guide the model’s inference.

Tables 9–11 reveal a task-dependent impact of Few-Shot prompting.

C.1.2 Impact of Chain-of-Thought Prompting

For Chain-of-Thought (CoT) prompting, we tailored task-specific prompts for each of the seven categories across all four tasks. The corresponding experimental results for GPT-5.1 are summarized in Table 12.

Overall, CoT prompting yields broad gains for the MLLM across categories, improving detection and localization performance on all except Manipulation Classification (MC). Averaged over seven categories, CoT increases Forgery Scope Discrimination (FSD) accuracy (+4.38%) and improves Textual Artifact Recognition (TAR) (+3.33%), while boosting Tampering Pinpointing (TP) localization quality (+4.25%), albeit with a slight rise in over-localization (+1.07%). These trends indicate that explicit reasoning steps mainly help the model surface more reliable authenticity cues and provide more correct localization hits, but may also encourage slightly more aggressive region proposals.

C.2 Analysis of Post-Processing Results

As shown in Table 13, to examine the robustness of models under common post-processing perturbations, we applied three representative image post-processing strategies. Specifically, we considered **Gaussian blurring**, **JPEG compression**, and **image scaling**. Gaussian blurring was implemented with a fixed blur radius of 5 to suppress high-frequency details. JPEG compression was applied with a quality factor of 50, introducing moderate compression artifacts. For image scaling, inputs were uniformly downsampled by a factor of 0.5 using the LANCZOS resampling method to preserve interpolation quality.

D Case Study

As shown in Figures 13–24, we provide 12 representative cases selected from the outputs of GPT-5.1. These cases span all seven academic image categories (Chart, Diagram, Micrograph, Stained Micrograph, Physical Object, Medical Imaging, and Others), and are systematically organized across all forensic dimensions.

E Use of AI Assistants

This research was driven entirely by the authors, who provided all core scientific insights, experimental designs, and analyses. We acknowledge the use of AI assistants during the preparation of this manuscript: Cursor was utilized to aid in code writing and data processing, large language models (LLMs) assisted with language editing to enhance readability, and generative models played a role in the construction of our synthetic data (Table 5). Furthermore, GPT-4o mini was employed to semantically reconstruct the captions of authentic academic images in *Text Constraint Fabrication*, with detailed methodologies provided in Appendix A.5.2. We emphasize that all AI-assisted content was thoroughly reviewed and validated by the authors, who bear full responsibility for the scientific integrity and fundamental contributions of this work.

Model	Count	Access Links
Text Constraint Fabrication		
Janus-Pro-7B	364	https://huggingface.co/deepseek-ai/Janus-Pro-7B
DALL-E 3	273	https://openai.com/zh-Hans-CN/index/dall-e-3/
FLUX.1 Dev	273	https://huggingface.co/black-forest-labs/FLUX.1-dev
FLUX.2 Pro	70	https://openrouter.ai/black-forest-labs/flux.2-pro
Midjourney V6	249	https://www.midjourney.com
Midjourney V7	247	https://www.midjourney.com
Stable Diffusion 3.5 Large Turbo	248	https://huggingface.co/stabilityai/stable-diffusion-3.5-large-turbo
Imagen 4	128	https://deepmind.google/models/imagen
Nano Banana Pro	71	https://deepmind.google/models/gemini-image/pro/
klimg-v2	190	https://klimgai.com/
CogView-4	183	https://docs.z.ai/guides/image/cogview-4
Jimeng 3.0	145	https://jimeng.jiaying.com/
Seedream 3.0	128	https://console.volcengine.com/ark/region:ark+cn-beijing/model/detail?Id=doubao-seedream-3-0-t2i
Wan2.1-T2I-Turbo	191	https://bailian.console.aliyun.com/cn-beijing?tab=model#/model-market/detail/wanx2.1-t2i-turbo?serviceSite=asia-pacific-china
Wan2.5-T2I-Preview	53	https://bailian.console.aliyun.com/cn-beijing?tab=model#/model-market/detail/wan2.5-t2i-preview?serviceSite=asia-pacific-china
Qwen-Image-Plus	35	https://bailian.console.aliyun.com/cn-beijing?tab=model#/model-market/detail/qwen-image-plus?serviceSite=asia-pacific-china
SenseNova v6 MiaoHua	273	https://www.sensenova.cn/en
Image Inference Forgery		
DALL-E 2	292	https://openai.com/zh-Hans-CN/index/dall-e-2/
FLUX.1 Schnell	482	https://huggingface.co/black-forest-labs/FLUX.1-schnell
Stable Diffusion 3.5 Medium	462	https://huggingface.co/stabilityai/stable-diffusion-3.5-medium
klimg-v1	253	https://klimgai.com/
Seedream 4.5	69	https://console.volcengine.com/ark/region:ark+cn-beijing/model/detail?Id=doubao-seedream-4-5
SenseNova v6 MiaoHua	506	https://www.sensenova.cn/en
Hunyuan 2.0	210	https://aistudio.tencent.com/modelSquare/home/list?modelKey=Multimodal
Targeted Region Restoration		
RePaint	1,650	https://github.com/andreas128/RePaint
Targeted Region Editing		
SenseNova v6 MiaoHua v2	701	https://www.sensenova.cn/en
Stable Diffusion	325	https://platform.stability.ai/docs/api-reference#tag/Edit/paths/~1v2beta~1stable-image~1edit~1inpaint/post
GPT Image 1	139	https://developers.openai.com/api/docs/models/gpt-image-1

Table 5: Generative models used for each forgery strategy.

Model	Capabilities	Access Links
Hybrid MLLM-Assisted Expert Models		
FakeShield	Detection + Localization	https://github.com/zhipeixu/FakeShield
SIDA	Detection + Localization	https://github.com/hzlsaber/SIDA
FakeVLM	Detection	https://github.com/opendatalab/FakeVLM
Vision-Only Expert Models		
DRCT	Detection	https://github.com/beibuwandeluori/DRCT
DIRE	Detection	https://github.com/ZhendongWang6/DIRE
AIDE	Detection	https://github.com/shilinyan99/AIDE
AlignedForensics	Detection	https://github.com/AniSundar18/AlignedForensics
Denoising Trajectory	Detection	https://github.com/Eleven4AI/DTAD
DDA	Detection	https://github.com/roy-ch/Dual-Data-Alignment

Table 6: Detailed information of benchmarked expert models.

Provider	Model	Size	Access Links
Proprietary MLLMs			
OpenAI	GPT-5.1	–	https://openrouter.ai/openai/gpt-5.1
	GPT-4.1	–	https://openrouter.ai/openai/gpt-4.1
	OpenAI o4-mini-high	–	https://openrouter.ai/openai/o4-mini-high
Anthropic	Claude Sonnet 4.5	–	https://openrouter.ai/anthropic/claude-sonnet-4.5
Google	Gemini 3 Pro Preview	–	https://openrouter.ai/google/gemini-3-pro-preview
	Gemini 2.5 Flash	–	https://openrouter.ai/google/gemini-2.5-flash
ByteDance	Doubao-Seed-1.6	–	https://console.volcengine.com/ark/model/detail?Id=doubao-seed-1-6
	Doubao-Seed-1.6-thinking	–	https://console.volcengine.com/ark/model/detail?Id=doubao-seed-1-6-thinking
	Doubao-Seed-1.6-flash	–	https://console.volcengine.com/ark/model/detail?Id=doubao-seed-1-6-flash
	Doubao-1.5-thinking-vision-pro	–	https://console.volcengine.com/ark/model/detail?Id=doubao-1-5-thinking-vision-pro
	Doubao-1.5-vision-pro	–	https://console.volcengine.com/ark/model/detail?Id=doubao-1-5-vision-pro
Alibaba	Qwen3-VL-Plus	–	https://bailian.console.aliyun.com/cn-beijing?tab=model#/model-market/detail/qwen3-vl-plus?serviceSite=asia-pacific-china
	Qwen-VL-Max	–	https://bailian.console.aliyun.com/cn-beijing?tab=model#/model-market/detail/qwen-vl-max?serviceSite=asia-pacific-china
	Qwen-VL-Plus	–	https://bailian.console.aliyun.com/cn-beijing?tab=model#/model-market/detail/qwen-vl-plus?serviceSite=asia-pacific-china
Open-Source MLLMs			
Google	Gemma 3 27B	27B	https://openrouter.ai/google/gemma-3-27b-it
Alibaba	Qwen2.5-VL-72B	72B	https://bailian.console.aliyun.com/cn-beijing?tab=model#/model-market/detail/qwen2.5-vl-72b-instruct?serviceSite=asia-pacific-china
Meta	Llama 4 Maverick	400B-A17B	https://openrouter.ai/meta-llama/llama-4-maverick
Mistral AI	Minstral 3 14B	14B	https://openrouter.ai/mistralai/minstral-14b-2512
LLaVA Community	LLaVA-NeXT-7B	7B	https://huggingface.co/llava-hf/llava-v1.6-vicuna-7b-hf
	LLaVA-NeXT-34B	34B	https://huggingface.co/llava-hf/llava-v1.6-34b-hf
	LLaVA-Interleave-7B	7B	https://huggingface.co/llava-hf/llava-interleave-qwen-7b-hf
	LLaVA-Interleave-7B-DPO	7B	https://huggingface.co/llava-hf/llava-interleave-qwen-7b-dpo-hf
	LLaVA-OneVision-7B	7B	https://huggingface.co/lmms-lab/llava-onevision-qwen2-7b-ov
	LLaVA-OneVision-72B	72B	https://huggingface.co/lmms-lab/llava-onevision-qwen2-72b-ov
	LLaVA-OneVision-7B-Chat	7B	https://huggingface.co/lmms-lab/llava-onevision-qwen2-7b-ov-chat
Unified Multimodal Understanding and Generation Model			
DeepSeek	Janus-Pro-7B	7B	https://huggingface.co/deepseek-ai/Janus-Pro-7B

Table 7: Detailed information of benchmarked MLLMs.

Model	$\gamma = 0.3$		$\gamma = 0.5$		$\gamma = 0.7$	
	NFI	Rank	NFI	Rank	NFI	Rank
Proprietary MLLMs						
GPT-5.1	50.09	1	48.80	1	47.54	1
GPT-4.1	43.32	5	43.31	5	43.30	5
OpenAI o4-mini-high	43.25	6	42.77	7	42.30	7
Claude Sonnet 4.5	28.37	18	26.83	19	25.38	19
Gemini 3 Pro Preview	47.89	3	45.79	4	43.78	4
Gemini 2.5 Flash	48.47	2	47.02	2	45.62	3
Doubao-Seed-1.6	41.75	9	41.73	9	41.70	9
Doubao-Seed-1.6-thinking	46.72	4	46.66	3	46.60	2
Doubao-Seed-1.6-flash	36.16	12	36.03	11	35.90	11
Doubao-1.5-thinking-vision-pro	42.54	8	42.39	8	42.25	8
Doubao-1.5-vision-pro	34.21	14	34.07	14	33.94	14
Qwen3-VL-Plus	16.55	22	16.53	21	16.51	21
Qwen-VL-Max	43.12	7	43.01	6	42.89	6
Qwen-VL-Plus	33.98	15	33.89	16	33.80	16
Open-Source MLLMs						
Gemma 3 27B	40.16	10	37.55	10	35.11	13
Qwen2.5-VL-72B	38.75	11	38.71	10	38.67	10
Llama 4 Maverick	31.01	16	29.15	17	27.40	17
Ministral 3 14B	21.51	20	19.90	20	18.40	20
LLaVA-NeXT-7B	19.54	21	16.53	21	13.99	24
LLaVA-NeXT-34B	28.37	18	27.74	18	27.12	18
LLaVA-Interleave-7B	13.81	23	13.61	24	13.41	24
LLaVA-Interleave-7B-DPO	28.61	17	28.18	16	27.75	16
LLaVA-OneVision-7B	11.63	26	11.45	26	11.28	26
LLaVA-OneVision-72B	35.10	13	33.81	15	32.57	15
LLaVA-OneVision-7B-Chat	12.36	25	12.19	25	12.03	25
Unified Multimodal Understanding and Generation Model						
Janus-Pro-7B	13.78	24	13.77	23	13.76	23

Table 8: Sensitivity of NFI to the OLR penalty exponent γ . Lower ranks indicate better performance.

Model	Prompting	ACC	M-F1	No-Forgery	Entire-F1	Partial-F1
GPT-5.1	None	50.99	46.16	52.54	74.70	11.25
	Few-Shot	64.67 (+13.68)	57.67 (+11.51)	59.34 (+6.80)	79.78 (+5.08)	33.90 (+22.65)
Doubao-Seed-1.6-thinking	None	44.57	37.13	41.97	66.44	2.97
	Few-Shot	62.60 (+18.03)	55.51 (+18.38)	57.06 (+15.09)	76.72 (+10.28)	32.75 (+29.78)
Qwen3-VL-Plus	None	38.77	35.72	41.23	50.12	15.82
	Few-Shot	69.97 (+31.20)	65.20 (+29.48)	50.42 (+9.19)	80.07 (+29.95)	65.09 (+49.27)

Table 9: **Forgery Scope Discrimination** performance under **Few-Shot prompting** strategies.

Model	Prompting	ACC	M-F1	Yes-F1	No-F1
GPT-5.1	None	76.43	76.38	77.50	75.25
	Few-Shot	84.23 (+7.80)	83.37 (+6.99)	87.16 (+9.66)	79.57 (+4.32)
Doubao-Seed-1.6-thinking	None	77.92	77.57	80.38	74.76
	Few-Shot	80.25 (+2.33)	78.98 (+1.41)	84.14 (+3.76)	73.83 (-0.93)
Qwen3-VL-Plus	None	79.25	78.98	81.38	76.59
	Few-Shot	81.24 (+1.99)	80.63 (+1.65)	84.06 (+2.68)	77.19 (+0.60)

Table 10: **Textual Artifact Recognition** performance under **Few-Shot prompting** strategies.

Model	Prompting	ACC	M-F1	Insert-F1	Remove-F1	Alter-F1
GPT-5.1	None	60.07	48.34	34.93	36.65	73.43
	Few-Shot	27.22 (-32.85)	30.53 (-17.81)	18.18 (-16.75)	42.25 (+5.60)	31.15 (-42.28)
Doubao-Seed-1.6-thinking	None	58.47	41.85	26.10	26.00	73.43
	Few-Shot	25.80 (-32.67)	23.75 (-18.10)	23.88 (-2.22)	23.46 (-2.54)	23.90 (-49.53)
Qwen3-VL-Plus	None	59.28	36.89	15.54	20.37	74.75
	Few-Shot	33.75 (-25.53)	31.52 (-5.37)	20.80 (+5.26)	25.60 (+5.23)	48.15 (-26.60)

Table 11: **Manipulation Classification** performance under **Few-Shot prompting** strategies.

Category	Prompting	Detection				Localization			
		FSD		TAR		MC		TP	
		ACC	M-F1	ACC	M-F1	ACC	M-F1	CLA	OLR
Chart	None	41.10	41.99	77.34	77.30	70.73	48.50	50.47	10.17
	CoT	50.50 (+9.40)	44.98 (+2.99)	80.81 (+3.47)	79.77 (+2.47)	62.42 (-8.31)	50.64 (+2.14)	53.67 (+3.20)	11.86 (+1.69)
Diagram	None	65.35	55.49	80.19	79.35	58.57	46.57	40.50	11.83
	CoT	73.35 (+8.00)	62.15 (+6.66)	83.75 (+3.56)	82.60 (+3.25)	47.31 (-11.26)	40.22 (-6.35)	50.81 (+10.31)	10.27 (-1.56)
Micrograph	None	39.09	43.33	71.98	64.39	55.30	45.05	46.13	13.18
	CoT	56.78 (+17.69)	52.02 (+8.69)	73.95 (+1.97)	68.85 (+4.46)	36.73 (-18.57)	35.03 (-10.02)	51.36 (+5.23)	15.18 (+2.00)
Stained Micrograph	None	34.79	42.16	57.72	52.90	50.47	47.21	53.43	16.31
	CoT	52.26 (+17.47)	48.68 (+6.52)	58.76 (+1.04)	55.83 (+2.93)	35.20 (-15.27)	34.46 (-12.75)	54.61 (+1.18)	10.64 (-5.67)
Physical Object	None	48.49	47.81	60.77	57.93	59.23	42.01	42.58	9.46
	CoT	57.66 (+9.17)	51.05 (+3.24)	65.74 (+4.97)	64.39 (+6.46)	40.16 (-19.07)	36.64 (-5.37)	40.65 (-1.93)	14.84 (+5.38)
Medical Imaging	None	48.92	50.31	75.12	74.78	59.47	45.47	49.81	13.69
	CoT	64.84 (+15.92)	60.42 (+10.11)	79.66 (+4.54)	79.64 (+4.86)	32.47 (-27.00)	32.60 (-12.87)	54.02 (+4.21)	21.84 (+8.15)
Others	None	61.89	51.55	83.00	82.78	66.67	57.37	48.40	13.24
	CoT	67.80 (+5.91)	59.61 (+8.06)	86.79 (+3.79)	86.57 (+3.79)	58.33 (-8.34)	55.46 (-1.91)	55.95 (+7.55)	10.71 (-2.53)
Avg.	None	48.52	47.52	72.30	69.92	60.06	47.45	47.33	12.55
	CoT	52.90 (+4.38)	47.36 (-0.16)	75.64 (+3.33)	73.95 (+4.03)	44.66 (-15.40)	40.72 (-6.73)	51.58 (+4.25)	13.62 (+1.07)

Table 12: Performance comparison of **GPT-5.1** across categories under **CoT** prompting strategies.

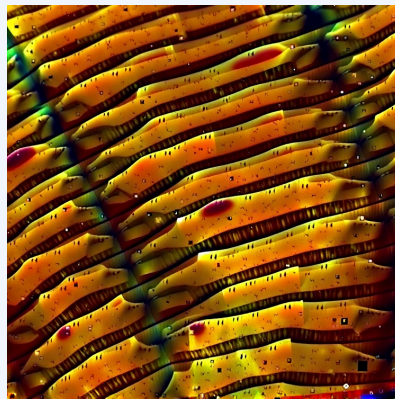
Model	Post-Processing	Detection				Localization					
		FSD		TAR		MC		TP			
		ACC	M-F1	ACC	M-F1	ACC	M-F1	Region-Level		Pixel-Level	
							CLA	OLR	IoU	F1	
GPT-5.1	None	50.99	46.16	76.43	76.38	60.07	48.34	46.87	12.25	-	-
	Gauss	24.22 (-26.77)	25.74 (-20.42)	46.32 (-30.11)	42.62 (-33.76)	58.74 (-1.33)	40.22 (-8.12)	21.58 (-25.29)	13.03 (+0.78)	-	-
	JPEG	47.22 (-3.77)	40.38 (-5.78)	72.16 (-4.27)	72.13 (-4.25)	48.99 (-11.08)	39.91 (-8.43)	9.39 (-37.48)	0.11 (-12.14)	-	-
	Scaling	49.49	41.25	71.00	70.99	59.40	46.91	11.53	0.12	-	-
		(-1.50)	(-4.91)	(-5.43)	(-5.39)	(-0.67)	(-1.43)	(-35.34)	(-12.13)	-	-
FakeShield	None	59.72	62.08	-	-	-	-	-	-	30.09	69.11
	Gauss	36.73	34.26	-	-	-	-	-	-	20.77	29.73
		(-22.99)	(-27.82)	-	-	-	-	-	-	(-9.32)	(-39.38)
	JPEG	40.26	36.15	-	-	-	-	-	-	22.07	28.16
		(-19.46)	(-25.93)	-	-	-	-	-	-	(-8.02)	(-40.95)
Scaling	38.01	35.04	-	-	-	-	-	-	22.82	29.09	
	(-21.71)	(-27.04)	-	-	-	-	-	-	(-7.27)	(-40.02)	

Table 13: **Impact of post-processing perturbations.**



Forgery Scope Discrimination

Category: Micrograph



Question

The task requires the model to act as an image forensics expert and identify the forgery scope of the given academic image panel, strictly following the specified output format.

Options

- A. Entire Forgery — entirely AI-generated
- B. Partial Forgery — partially AI-edited (with localized forgery)
- C. No Forgery
- D. Not Sure — too uncertain to judge

Issue A. Entire Forgery



(GPT-5.1) Model Answer: A. Entire Forgery

Model Explanation

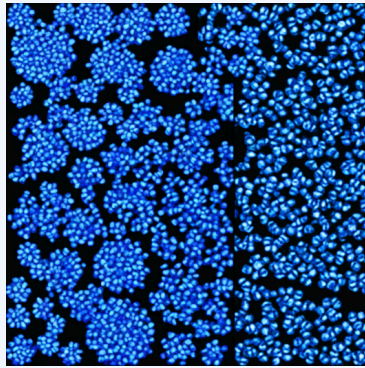
The entire image exhibits uniform AI-art-like textures, exaggerated colors, and repetitive artifacts inconsistent with natural scientific imagery, indicating it is fully AI-generated rather than partially edited.

Figure 13: Success case by GPT-5.1 on Forgery Scope Discrimination. The panel belongs to the Micrograph category.



Forgery Scope Discrimination

Category: Stained Micrograph



Question

The task requires the model to act as an image forensics expert and identify the forgery scope of the given academic image panel, strictly following the specified output format.

Options

- A. Entire Forgery — entirely AI-generated
- B. Partial Forgery — partially AI-edited (with localized forgery)
- C. No Forgery
- D. Not Sure — too uncertain to judge

Issue A. Entire Forgery



(GPT-5.1) Model Answer: D. Not Sure **X**

Model Explanation

The image shows plausible microscopic structures without obvious local editing artifacts or clear generative irregularities, but there are also no definitive cues to confirm or rule out full AI generation, making the level of AI intervention uncertain.



Analysis

Error Analysis

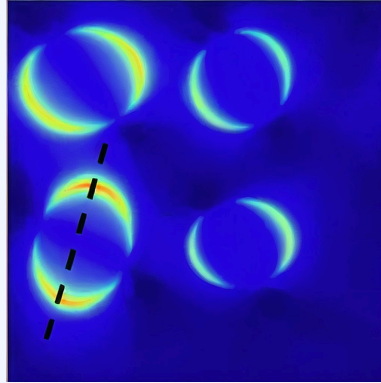
In the absence of explicit local editing artifacts, models tend to adopt a conservative abstention strategy, treating visually plausible and globally consistent structures as insufficient evidence for full-image synthesis. This indicates an over-reliance on localized forensic cues and an underutilization of global distributional signals, causing high-realism global forgeries to be underestimated when salient manipulation triggers are absent.

Figure 14: **Failure case by GPT-5.1 on Forgery Scope Discrimination.** The panel belongs to the Stained Micrograph category.



Forgery Scope Discrimination

Category: Diagram



Question

The task requires the model to act as an image forensics expert and identify the forgery scope of the given academic image panel, strictly following the specified output format.

Options

- A. Entire Forgery — entirely AI-generated
- B. Partial Forgery — partially AI-edited (with localized forgery)
- C. No Forgery
- D. Not Sure — too uncertain to judge

Issue A. Entire Forgery



(GPT-5.1) Model Answer: C. No Forgery **X**

Model Explanation

The image appears to be a standard scientific heatmap/visualization with smooth gradients and clean geometric shapes, without visible artifacts or inconsistencies suggestive of AI-generated or edited regions.



Analysis

Error Analysis

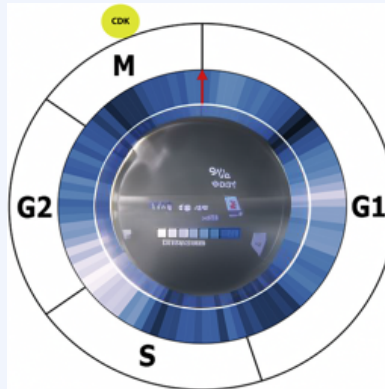
The model interprets the smooth gradients, clean geometric structures, and overall visual coherence as evidence of authenticity, conflating stylistic regularity with physical plausibility. This reflects a bias toward treating well-structured heatmap-like renderings as inherently trustworthy, while overlooking global generative cues such as overly idealized symmetry, uniform smoothness, and the absence of acquisition-related noise.

Figure 15: Failure case by GPT-5.1 on Forgery Scope Discrimination. The panel belongs to the Diagram category.



Forgery Scope Discrimination

Category: Chart



Question

The task requires the model to act as an image forensics expert and identify the forgery scope of the given academic image panel, strictly following the specified output format.

Options

- A. Entire Forgery — entirely AI-generated
- B. Partial Forgery — partially AI-edited (with localized forgery)
- C. No Forgery
- D. Not Sure — too uncertain to judge

Issue B. Partial Forgery



(GPT-5.1) Model Answer: C. No Forgery **X**

Model Explanation

The chart and embedded image elements appear consistent with standard scientific figure design, with clean boundaries, regular text, and no obvious AI-generated artifacts or irregular patterns.



Analysis

Error Analysis

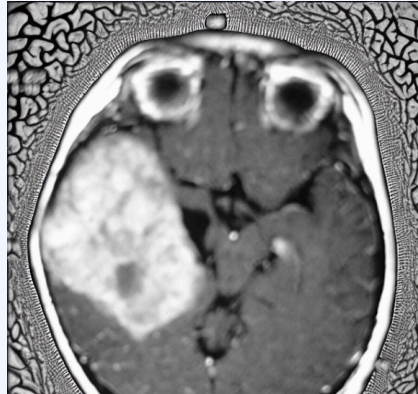
The model interprets the clean layout, regular annotations, and well-aligned geometric components as indicators of authenticity, overlooking subtle inconsistencies introduced by localized manipulation. This reflects an over-reliance on global visual coherence and design regularity, causing localized forgeries embedded within otherwise well-structured scientific figures to be missed.

Figure 16: Failure case by GPT-5.1 on Forgery Scope Discrimination. The panel belongs to the Chart category.



Forgery Scope Discrimination

Category: Medical Imaging



Question

The task requires the model to act as an image forensics expert and identify the forgery scope of the given academic image panel, strictly following the specified output format.

Options

- A. Entire Forgery — entirely AI-generated
- B. Partial Forgery — partially AI-edited (with localized forgery)
- C. No Forgery
- D. Not Sure — too uncertain to judge

Issue B. Partial Forgery



(GPT-5.1) Model Answer: A. Entire Forgery **X**

Model Explanation

The entire image shows uniform synthetic textures, anatomically inconsistent brain structures, and repetitive decorative patterns around the border, indicating it is fully AI-generated rather than a localized edit of a real scan.



Analysis

Error Analysis

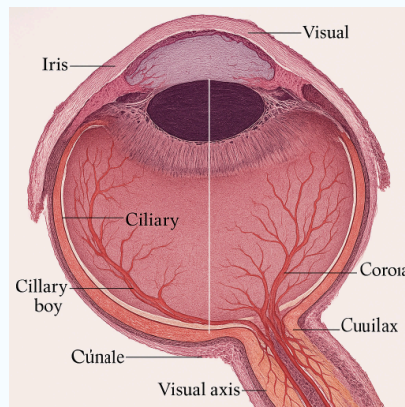
In this case, the model overemphasizes global visual irregularities—such as atypical textures, perceived anatomical inconsistencies, and decorative border patterns—and incorrectly attributes them to full-image synthesis. This leads to a misinterpretation of localized manipulations as evidence of global generation, blurring the distinction between partial edits and entirely synthetic content.

Figure 17: Failure case by GPT-5.1 on Forgery Scope Discrimination. The panel belongs to the Medical Imaging category.



Textual Artifact Recognition

Category: Diagram



Question

The task requires the model to assess whether any embedded textual elements in the given academic image exhibit signs of AI generation or AI editing, based solely on visual artifacts, structural irregularities, or semantic mismatches with the surrounding scientific content. The model should return a binary judgment (Yes/No), strictly following the specified output format.

Options

- A. Yes
- B. No

Issue A. Yes



(GPT-5.1) Model Answer: A. Yes

Model Explanation

Several labels (e.g., for anatomical parts of the eye) are misspelled or use nonsensical terms that do not correspond to established ocular anatomy, indicating semantic inconsistencies typical of AI-generated or altered text. The typography and label placement appear regular, but the incorrect terminology suggests the text is not authentic scientific labeling.

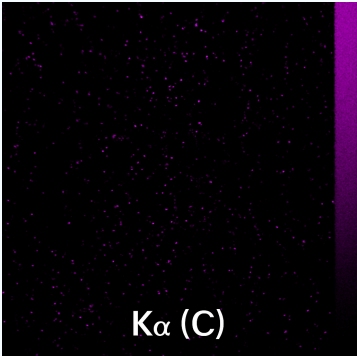
Figure 18: Success case by GPT-5.1 on Textual Artifact Recognition. The panel belongs to the Diagram category.

✘

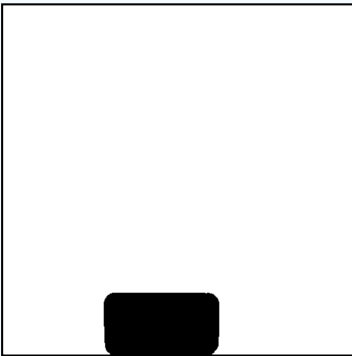
Textual Artifact Recognition

Category: Stained Micrograph

Input Image



GT Mask




Question

The task requires the model to assess whether any embedded textual elements in the given academic image exhibit signs of AI generation or AI editing, based solely on visual artifacts, structural irregularities, or semantic mismatches with the surrounding scientific content. The model should return a binary judgment (Yes/No), strictly following the specified output format.

Options


A. Yes
B. No

Issue A. Yes

 **(GPT-5.1) Model Answer: B. No** ✘

Model Explanation

The embedded text appears clean, well-aligned, and stylistically consistent with standard scientific figure labeling, with no visual artifacts or semantic mismatches relative to the elemental map image content.

 **Analysis**

Error Analysis

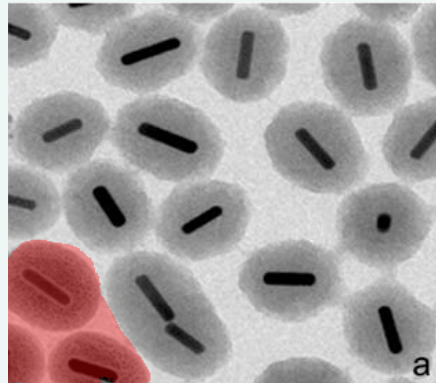
The failure stems from an over-reliance on semantic plausibility and visual neatness of the label (“Kα (C)”), while neglecting localized text-region artifacts. Specifically, the embedded text exhibits abnormal contrast and blending behavior against the dark elemental map background, with unnatural edge smoothness and luminance consistency that deviate from typical staining or elemental mapping annotations. The model’s explanation focuses on global alignment and stylistic cleanliness, overlooking text–background inconsistency cues that are critical for textual-level forgery detection in scientific micrographs.

Figure 19: **Failure case by GPT-5.1 on Textual Artifact Recognition.** The panel belongs to the Stained Micrograph category and contains a localized forgery generated via *Targeted Region Editing*.



Manipulation Classification

Category: Micrograph



Question

The task requires the model to identify the manipulation type of an AI-edited local region based on observable visual changes and structural inconsistencies, classifying the modification as Remove, Insert, or Alter, strictly following the specified output format.

Original Description

...high-magnification electron microscope image showing what appears to be bacterial cells or possibly microscopic particles. The image displays multiple oval or round-shaped structures with light gray borders and darker gray interiors. Within each of these oval structures is a distinct dark rod-shaped element...

Options

- A. Insert
- B. Remove
- C. Alter
- D. Not Sure

Issue **C. Alter**



(GPT-5.1) Model Answer: C. Alter

Model Explanation

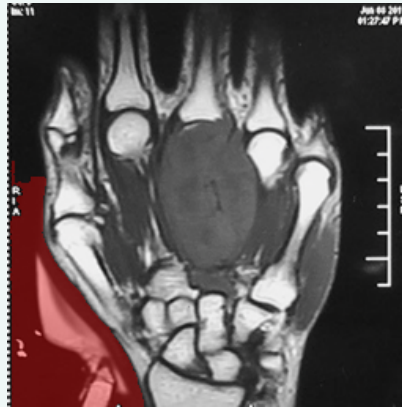
The red-highlighted area still contains similar cell-like structures but with altered shape/contrast compared with the surrounding ones, indicating modification of existing objects rather than adding new ones or removing them.

Figure 20: **Success case by GPT-5.1 on Manipulation Classification.** The panel belongs to the Micrograph category and contains a localized forgery generated via *Targeted Region Editing*.



Manipulation Classification

Category: Medical Imaging



Question

The task requires the model to identify the manipulation type of an AI-edited local region based on observable visual changes and structural inconsistencies, classifying the modification as Remove, Insert, or Alter, strictly following the specified output format.

Original Description

...MRI scan of a hand in coronal (frontal) view. The scan clearly displays the bony structures of the hand including multiple metacarpal bones, phalanges, and carpal bones at the wrist. There is a prominent dark (hypointense) mass visible in the central palm region, likely representing a soft tissue mass or cyst. The mass appears well-circumscribed and is pushing adjacent structures aside...

Options

- A. Insert
- B. Remove
- C. Alter
- D. Not Sure

Issue **A. Insert**



(GPT-5.1) Model Answer: **B. Remove** ❌

Model Explanation

The red-highlighted region overlays and occludes part of the original MRI anatomy, effectively removing existing image content rather than adding or modifying anatomical structures.



Analysis

Error Analysis

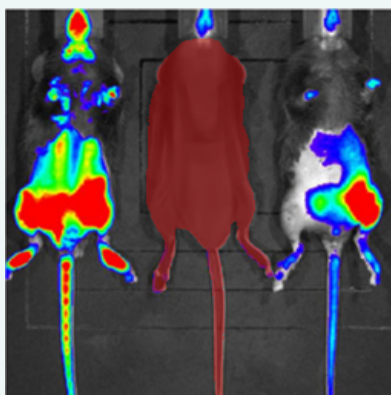
This error arises from a misinterpretation of occlusion as deletion, where the model assumes that the red-highlighted region merely covers or removes existing anatomical structures. In reality, the forged region introduces a new soft-tissue-like mass that was not present in the original image, thereby constituting an insertion rather than a removal. The model's explanation focuses on the apparent obstruction of underlying anatomy, overlooking the novel semantic content and abnormal anatomical presence introduced by the edited region. This case highlights the model's difficulty in distinguishing inserted foreground structures from removals in medical imaging scenarios.

Figure 21: Failure case by GPT-5.1 on Manipulation Classification. The panel belongs to the Medical Imaging category and contains a localized forgery generated via *Targeted Region Restoration*.



Manipulation Classification

Category: Medical Imaging



Question

The task requires the model to identify the manipulation type of an AI-edited local region based on observable visual changes and structural inconsistencies, classifying the modification as Remove, Insert, or Alter, strictly following the specified output format.

Original Description

...fluorescence imaging study involving three murine subjects, aligned in a horizontal arrangement. Each specimen exhibits varying intensity of fluorescent signals, depicted through a color spectrum ranging from red (indicating high signal intensity) to blue (indicating low signal intensity), the leftmost mouse displays pronounced red fluorescence across the abdomen and proximal limbs, suggesting a significant accumulation of the fluorescent marker...

Options

- A. Insert
- B. Remove
- C. Alter
- D. Not Sure

Issue **B. Remove**



(GPT-5.1) Model Answer: **A. Alter** ❌

Model Explanation

The central mouse, originally present as a foreground subject with distinct fluorescence patterns, has been visually modified and overpainted in the highlighted region, changing its appearance rather than removing it or adding a new object.



Analysis

Error Analysis

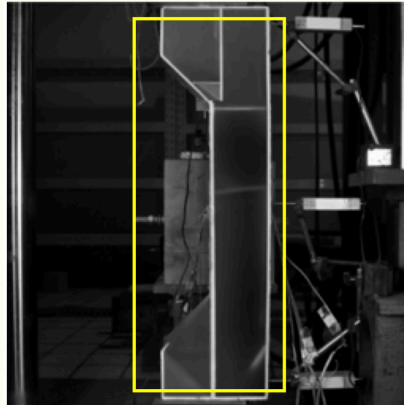
This error stems from confusing content erasure with appearance modification. In the red-highlighted region, the central mouse subject has been entirely eliminated, leaving a visually plausible but semantically empty or background-like region, which constitutes a removal rather than an alteration. However, the model interprets the uniform overlay and smooth intensity transition as evidence of repainting or overpainting, leading it to infer an alteration of existing content. This case reveals the model's tendency to over-attribute surface-level visual continuity to modification, while failing to recognize the absence of an expected foreground object.

Figure 22: Failure case by GPT-5.1 on Manipulation Classification. The panel belongs to the Medical Imaging category and contains a localized forgery generated via *Targeted Region Restoration*.



Tampering Pinpointing

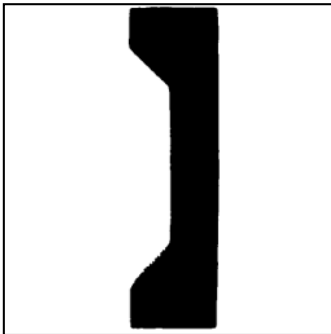
Category: Physical Object



Question

The task requires the model to precisely localize all AI-edited regions in the given image by outputting bounding boxes that tightly enclose each forged area, strictly following the specified output format.

Forged Region (GT Mask)



(GPT-5.1) Model Answer: `[[80, 10, 175, 245]]`

The predicted region is highlighted with the **yellow bounding box** in the image.

Model Explanation

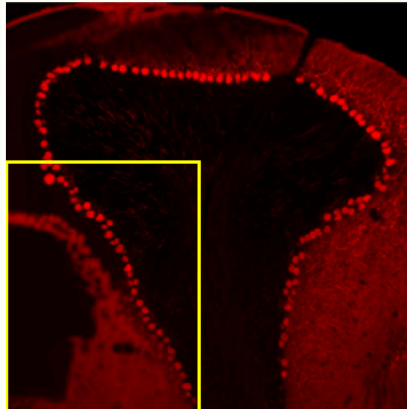
The central tall, glowing structural shape has unnaturally sharp uniform borders and smooth interior gradients that differ from the surrounding photographic background, indicating it was digitally inserted or modified.

Figure 23: **Success case by GPT-5.1 on Tampering Pinpointing.** The panel belongs to the Physical Object category and contains a localized forgery generated via *Targeted Region Restoration*.



Tampering Pinpointing

Category: Stained Micrograph



Question

The task requires the model to precisely localize all AI-edited regions in the given image by outputting bounding boxes that tightly enclose each forged area, strictly following the specified output format.

Forged Region (GT Mask)



(GPT-5.1) Model Answer: `[[0, 100, 120, 255]]`



The predicted region is highlighted with the **yellow bounding box** in the image.

Model Explanation

The lower-left region has an unnaturally sharp rectangular boundary and a flat dark fill that contrasts with the organic curved structures and textured signal of the rest of the image, indicating likely AI editing or compositing.



Analysis

Error Analysis

The model fails to accurately localize the forged region, producing a single coarse bounding box that only partially overlaps with the ground-truth mask. This error is primarily caused by an overly conservative and axis-aligned localization strategy, where the model prioritizes a rectangular region with high visual contrast rather than tightly following the irregular forged area. Although the model correctly detects the presence of abnormal visual cues (e.g., sharp boundaries and flat dark fills), it lacks the ability to adapt bounding box geometry to complex biological structures, resulting in under-coverage of curved and branching regions. This case highlights a key limitation of current MLLMs in fine-grained spatial grounding for non-rigid anatomical patterns commonly found in stained micrographs.

Figure 24: Failure case by GPT-5.1 on Tampering Pinpointing. The panel belongs to the Stained Micrograph category and contains a localized forgery generated via *Targeted Region Restoration*.