



Towards Zero-Shot Anomaly Detection and Reasoning with Multimodal Large Language Models

CVPR 2025 (Highlight)











Jiacong Xu

Shao-Yuan Lo

Bardia Safaei

Vishal M. Patel

Isht Dwivedi

4/12/2025

What is Anomaly Detection?

- **Problem definition**: An AD model is exclusively trained with **normal** data and is asked to identify whether a query example is **normal** or **anomalous**.
- Motivation: Anomalies are often rare and long-tailed, so they are costly to collect.
- Example:
 - Normal data: Flawless objects
 - Anomalous data: Defects





Conventional Learning-based AD Approaches

- Full-shot training: An AD model is trained by a large amount of normal data to learn normal patterns
- Output format: Anomaly scores -> Thresholding
- Metrics: AUROC (area under ROC curve)



The Era of Multimodal Large Language Models



How can AD benefit from MLLMs?

What breakthroughs can MLLMs bring to this long-established vision problem?

Remain underexplored!

Anomaly Detection in the Era of MLLMs



Towards Zero-Shot Anomaly Detection and Reasoning with Multimodal Large Language Models

Jiacong Xu¹* Shao-Yuan Lo² Bardia Safaei¹ Vishal M. Patel¹ Isht Dwivedi² ¹Johns Hopkins University ²Honda Research Institute USA {jxu155, bsafaei1, vpatel36}@jhu.edu {shao-yuan_lo, idwivedi}@honda-ri.com CVPR 2025 (Highlight)

- One of the **first zero-shot and reasoning** methods for IAD
 - => Explain why normal/anomaly
- One of the first instruction tuning datasets and benchmarks for IAD reasoning
 - => Allow us to train and evaluate models for their AD reasoning capabilities

Problem Statement

Anomaly detection

PCB Example: Excessive Solder

User

	UCC Trig Echo Grd
Are there any anomali	es visible in the image?

GPT-4o	The image shows an HC-SR04 ultrasonic sensor. There don't appear to be any visible anomalies with the sensor itself. The components and pins seem intact and properly aligned.
LLaVA-OV	The image shows an ultrasonic sensor module, specifi- cally the HC-SR04 model. There are no visible anoma- lies in the image; it appears to be a standard represen- tation of this type of sensor module.
Ours	Yes, there is a noticeable anomaly in the image: a short

circuit on the PCB, indicated by a white, crusty residue near the 'Echo' pin.

Table 3. Anomaly-OV presents more accurate anomaly detection.

Basic reasoning

Macaroni Example: Yellowish Spot

User	Do you see any anomalies in the image?
GPT-40	The image shows four pieces of elbow macaroni on a green background. The anomaly is that the macaroni pieces are not whole; they are cut in half.
LLaVA-OV	The image shows four pieces of pasta, specifically macaroni shells, arranged on a green textured surface. The pasta appears to be uniformly colored and shaped, with no visible defects or anomalies.
Ours	Yes, there is an anomaly in the image. The bottom right pasta piece has a noticeable yellowish discoloration on its surface.
able 5 Anor	ns surface.

Table 5. Anomaly-OV presents more precise anomaly reasoning.

Complex reasoning

In-the-Wild Example: Graffitied Road Sign



User	Describe the anomaly in the image.						
Anomaly-OV	The anomaly in the image is the graffiti on the high- way sign, which includes black spray paint marks and scribbles over the word 'Seattle' and other parts of the sign.						
User	What is the reason for the anomaly?						
Anomaly-OV	The graffiti is likely the result of vandalism, where indi- viduals have used spray paint to deface public property for personal expression or as a form of protest.						
User	How does the anomaly affect the functionality of the object?						
Anomaly-OV	The graffiti affects the sign's functionality by poten- tially obstructing the visibility of important informa- tion, such as the word 'Seattle,' which could confuse drivers and lead to navigational errors.						
User	How to avoid the anomaly in the future?						
Anomaly-OV	To avoid such anomalies, regular maintenance and cleaning of signs should be implemented. Addition- ally, community awareness programs can help reduce vandalism by promoting respect for public property.						

Method

- Two-step inspection for improving detection
- Simulate human's visual anomaly inspection



Method

- MLLM backbone: Visual encoder + tokenizer + projector + LLM
- LTFM: (Step 1) Find suspicious areas
- VT Selector: (Step 2) Confirm anomalies
- **Training stage 1**: Train the anomaly expert (LTFM + VT Selector)
- Training stage 2: Visual instruction tuning (projector + LLM)





Instruction Tuning Dataset

- Visual instruction tuning (VIT) for improving reasoning
- An automatic pipeline for constructing AD VIT datasets
- WebAD: 72k, 380 classes
- MVTec AD [Bergmann et al. CVPR'19] : 5k, 15 classes



Instruction Tuning Dataset

• WebAD: 72k, 380 classes



Anomalous

Normal

Instruction Tuning Dataset

Anomaly-Instruct-125k

• Four types of images



Q: Do you notice any anomalies in the image? A: Yes, the anomaly in the image is a fine, dark crack visible on the ceramic cup, starting from the rim and extending downward.

Q: How can I prevent similar anomalies in the future? A: Avoid rapid temperature changes, handle the cup carefully to minimize impacts, and consider using cups made with thicker, more durable materials.



Q: Are there any anomalies in this cable image? A: Yes, the green wire is missing its internal copper conductors, unlike the blue and brown wires, which contain complete sets of bundled copper strands.

Q: What could be the reason behind this anomaly? A: The absence of copper in the green wire could be due to manufacturing defects such as failure to insert the strands or a mechanical issue where the strands were removed or misplaced during production or handling.

Q: What does the image show? A: The image shows a densely cellular tissue sample with features indicating potential metastatic breast cancer.

Q: What suggests that there might be malignant activity in the tissue sample? A: Malignant activity is suggested by irregularly shaped cells, hyperchromatic nuclei, and an increased nuclear-to-cytoplasmic ratio, indicative of abnormal and increased DNA content typical in cancerous cells.



Q: What kind of object are you examining in the images?
A: I'm examining a round cookie with a golden-brown surface and visible chocolate chips embedded across its surface.



Q: Are there any significant anomalies detected on this cookie?

A: Yes, there are two significant anomalies on the cookie: a central defect

featuring a noticeable hole and an edge defect showing a gap or crack.

WebAD-72k

[Collected by us via Google Image Search] MVTec AD [Bergmann et al. CVPR'19]

BMAD [Bao et al. CVPR'24]



Reasoning Benchmark

- VisA-D&R (761 normal + 1000 anomaly samples)
- **Detection metrics**: Accuracy/precision/recall/F1-score
- Reasoning metrics: ROUGE-L, Sentence-BERT, GPT-score
- Q1 & Q2: Basic reasoning
- Q3 & Q4: Complex reasoning

```
Detection:
Q: Are there any defects for the object in the image?
Please reply with 'Yes' or 'No'.
Reasoning:
Q1: Do you observe any anomalies in the image?
Q2: Can you describe the anomalies you observed?
Q3: What is the potential cause for the anomalies?
Q4: How can such anomalies be prevented in the future?
```

Detection Results of Zero-Shot IAD Approaches

- MLLM backbone: LLaVA-OV [Li et al. 2024]
- Metric: AUROC

Model	Industrial Defects						Medical Anomalies			Average
	MVTec AD	VisA	AITEX	ELPV	BTAD	MPDD	BrainMRI	HeadCT	Br35H	i i ciugo
CLIP [73]	74.1	66.4	71.0	59.2	34.5	54.3	73.9	56.5	78.4	63.1
CoOp [108]	88.8	62.8	66.2	73.0	66.8	55.1	61.3	78.4	86.0	70.9
WinCLIP [38]	91.8	78.8	73.0	74.0	68.2	63.6	92.6	90.0	80.5	79.2
APRIL-GAN [11]	86.2	78.0	57.6	65.5	73.6	73.0	89.3	89.1	93.1	78.4
AnoVL [19]	<u>92.5</u>	79.2	<u>72.5</u>	70.6	80.3	68.9	88.7	81.6	88.4	80.3
AnomalyCLIP [110]	91.5	82.1	62.2	<u>81.5</u>	88.3	77.0	90.3	<u>93.4</u>	94.6	84.5
AdaCLIP [6]	89.2	<u>85.8</u>	64.5	79.7	<u>88.6</u>	76.0	94.8	91.4	97. 7	<u>85.3</u>
Ours	94.0	91.1	72.0	83.0	89.0	81.7	<u>93.9</u>	97.6	<u>95.5</u>	88.6

Detection and Reasoning Results of MLLMs

- MLLM backbone: LLaVA-OV
- LLaVA-OV-0.5B*: Fine-tuned on our Anomaly-Instruct-125k

Model	Anomaly Detection				Low-	level Reas	Complex Reasoning		
	Accuracy	Precision	Recall	F1-score	ROUGE-L	SBERT	GPT-Score	SBERT	GPT-Score
GPT-4V [71]	0.68	0.90	0.49	0.55	0.16	0.65	3.31	0.77	5.64
GPT-40 [72]	0.70	0.83	0.71	0.68	0.24	0.71	4.84	0.81	6.89
Qwen2-VL-2B [87]	0.65	0.87	0.55	0.59	0.22	0.55	1.94	0.74	4.26
Qwen2-VL-7B [87]	0.76	<u>0.91</u>	0.69	0.75	0.25	0.61	3.09	0.68	4.62
InternVL-2-8B [13]	0.74	0.78	0.81	0.76	0.23	0.73	3.69	0.80	5.08
InternVL-2-26B [13]	0.73	0.86	0.66	0.68	0.21	0.74	4.13	0.80	5.49
IXC-2.5-7B [101]	0.72	0.88	0.63	0.67	0.21	0.58	2.45	0.77	5.14
LLaVA-OV-0.5B [44]	0.54	0.70	0.19	0.28	0.20	0.63	2.54	0.81	4.34
LLaVA-OV-7B [44]	0.71	0.95	0.56	0.63	0.24	0.66	3.57	0.79	5.44
LLaVA-OV-0.5B*	0.71	0.77	0.84	0.76	0.31	0.70	3.69	0.82	5.31
Anomaly-OV-0.5B	0.79	0.86	0.83	0.82	0.33	0.72	3.87	<u>0.83</u>	5.67
Anomaly-OV-7B	0.79	0.83	0.86	0.83	0.34	<u>0.73</u>	<u>4.26</u>	0.84	<u>6.34</u>

Towards Zero-Shot Anomaly Detection and Reasoning with Multimodal Large Language Models

Jiacong Xu¹* Shao-Yuan Lo² Bardia Safaei¹ Vishal M. Patel¹ Isht Dwivedi² ¹Johns Hopkins University ²Honda Research Institute USA {jxu155, bsafaei1, vpatel36}@jhu.edu {shao-yuan_lo, idwivedi}@honda-ri.com CVPR 2025 (Highlight)

- One of the **first zero-shot and reasoning** methods for IAD
 - => Explain why normal/anomaly
- One of the first instruction tuning datasets and benchmarks for IAD reasoning
 - => Allow us to train and evaluate models for their AD reasoning capabilities