# Anomaly Detection in the Era of Multimodal Large Language Models

Shao-Yuan Lo

Research Scientist @ Honda Research Institute USA
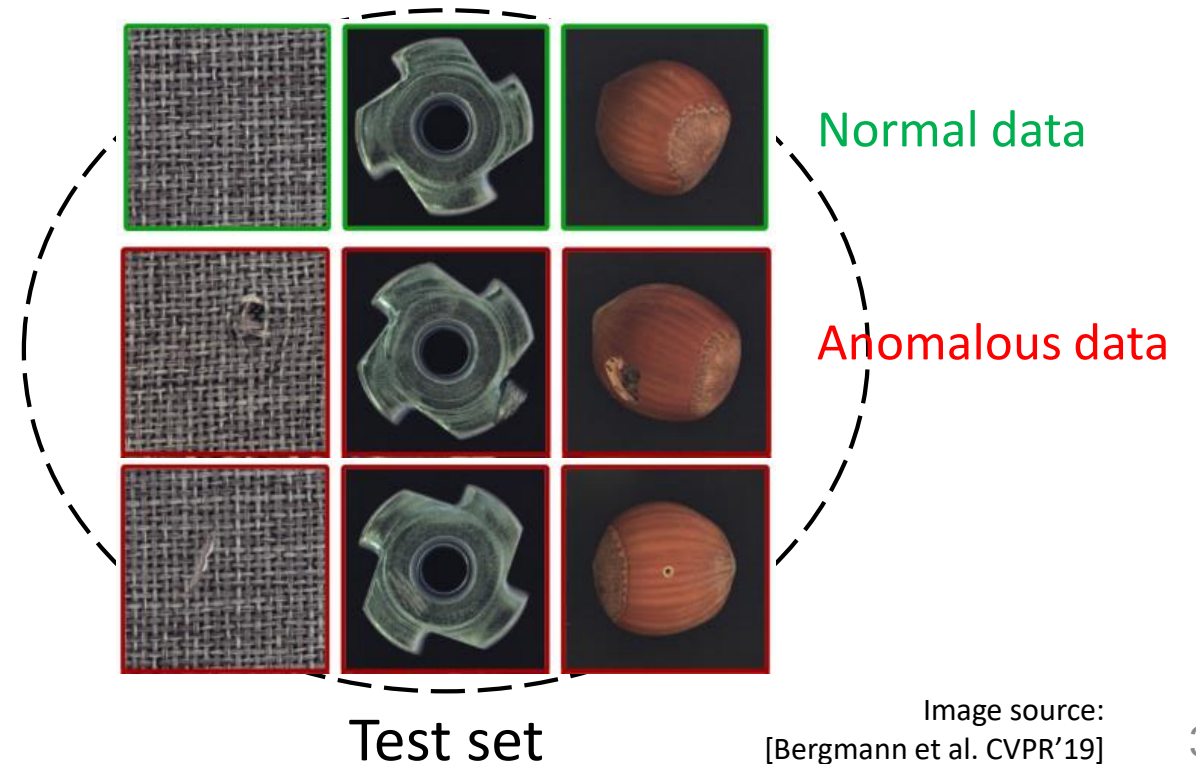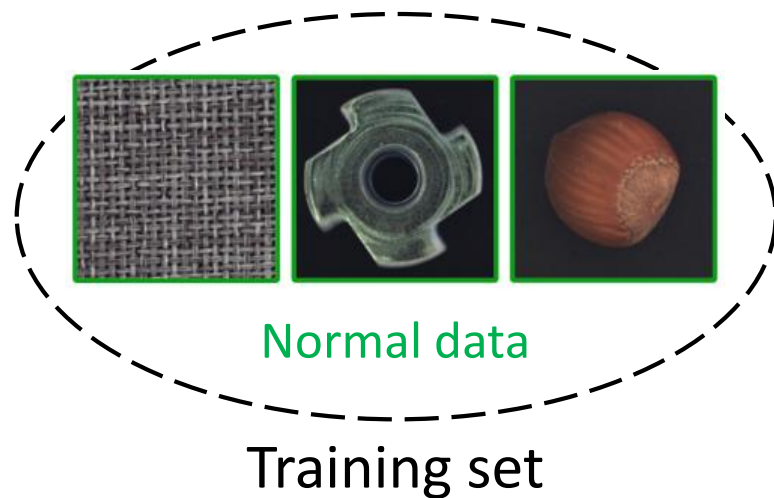
3/10/2025 @ UC Riverside

# About Me

- Research Scientist @ **Honda Research Institute USA**
  San Jose, CA (2023 - Present)

- Research Intern @ **Amazon**
  Seattle, WA (Summer 2021 & 2022)

- PhD in ECE @ **Johns Hopkins University**
  Baltimore, MD (2019 - 2023)

- MS in EE @ **National Chiao Tung University**
  Taiwan (2017 - 2019)

- BS in EECS @ **National Chiao Tung University**
  Taiwan (2013 - 2017)

# 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



Normal data

Training set

Normal data

Anomalous data

Test set

Image source:
[Bergmann et al. CVPR'19]
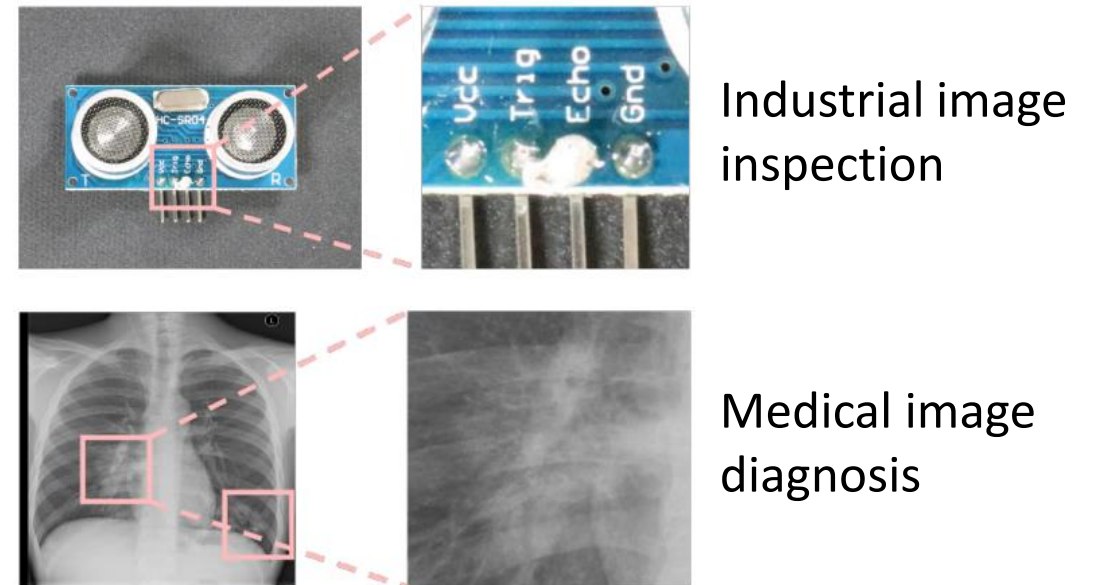
3

# Visual Anomaly Detection: Images and Videos

- Visual AD is a long-established problem in computer vision.
- Given its practical significance, AD has been widely deployed in various applications.
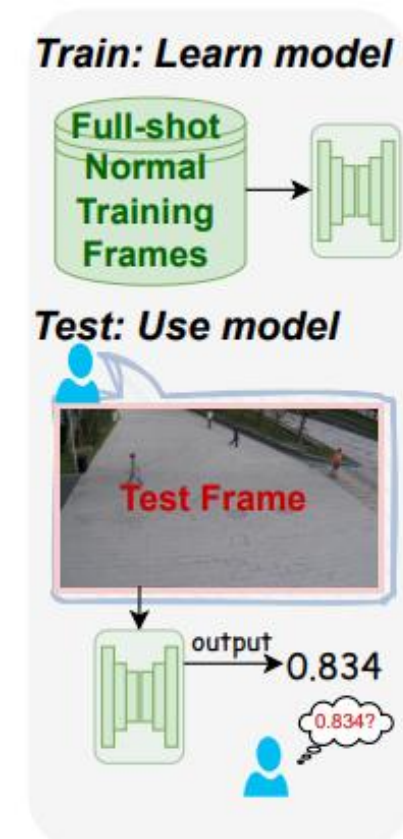


**Video Anomaly Detection (VAD)**

Security surveillance



**Image Anomaly Detection (IAD)**

Industrial image inspection

Medical image diagnosis

Image source: [Bergmann et al. CVPR'19] [Bao et al. CVPR'24] [Liu et al. CVPR'18]

# 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)



Train: Learn model

Full-shot Normal Training Frames

Test: Use model

Test Frame

output 0.834

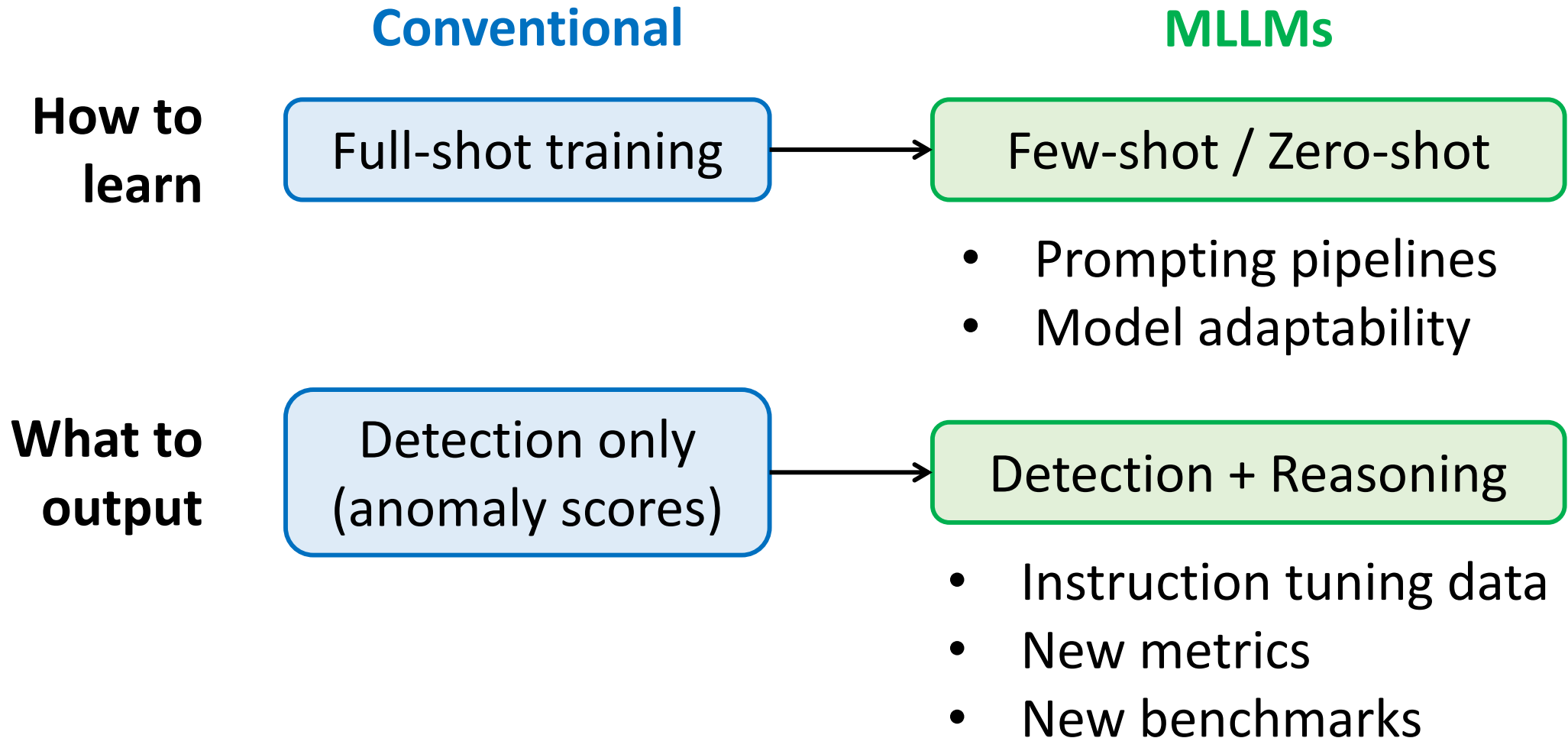0.834?

# 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

**Conventional**                    **MLLMs**

**How to learn**

Full-shot training → Few-shot / Zero-shot

- Prompting pipelines
- Model adaptability

**What to output**

Detection only (anomaly scores) → Detection + Reasoning

- Instruction tuning data
- New metrics
- New benchmarks

# Anomaly Detection in the Era of MLLMs

ECCV 2024

CVPR 2025

Video Anomaly Detection (VAD)

Image Anomaly Detection (IAD)

Security surveillance

Industrial image inspection

Medical image diagnosis

Image source: [Bergmann et al. CVPR'19] [Bao et al. CVPR'24] [Liu et al. CVPR'18]

# Follow the Rules: Reasoning for Video Anomaly Detection with Large Language Models

Yuchen Yang[1*], Kwonjoon Lee[2], Behzad Dariush[2], Yinzhi Cao[1], and Shao-Yuan Lo[2]

[1] Johns Hopkins University
{yc.yang, yinzhi.cao}@jhu.edu
[2] Honda Research Institute USA
{kwonjoon_lee, bdariush, shao-yuan_lo}@honda-ri.com

ECCV 2024

- One of the **first reasoning** methods for VAD
  - => Explain why normal/anomaly

- One of the **first few-shot prompting** methods for VAD
  - => Fast adaption to different definitions of "anomaly" for different applications

# Problem Statement

- **Assumption**: We only have a few <span style="color:green">**normal**</span> data for our specific application, and it's costly to collect <span style="color:red">**anomaly**</span> data.

- **Challenge**: The definition of "<span style="color:red">**anomaly**</span>" **depends on different context and downstream applications**.

- **Goal**: Develop a VAD model for our specific application (specific definition of "<span style="color:green">**normal**</span>" & "<span style="color:red">**anomaly**</span>") and explain the detection results.



person jogging **versus** person running outside a bank.

Image source: [Radu Tudor Ionescu. CVPR'23 Tutorial]

# Conventional Learning-based Approaches

- **Full-shot training**: A VAD model is trained by a large amount of normal data to learning normal patterns

- **Output format**: Anomaly scores -> Thresholding

- **Metrics**: AUROC (area under ROC curve)



*Train: Learn model*

Full-shot Normal Training Frames

*Test: Use model*

Test Frame

output 0.834

0.834?

# Query MLLMs Directly

- The implicit knowledge pre-trained in MLLMs may not align with specific VAD needs.

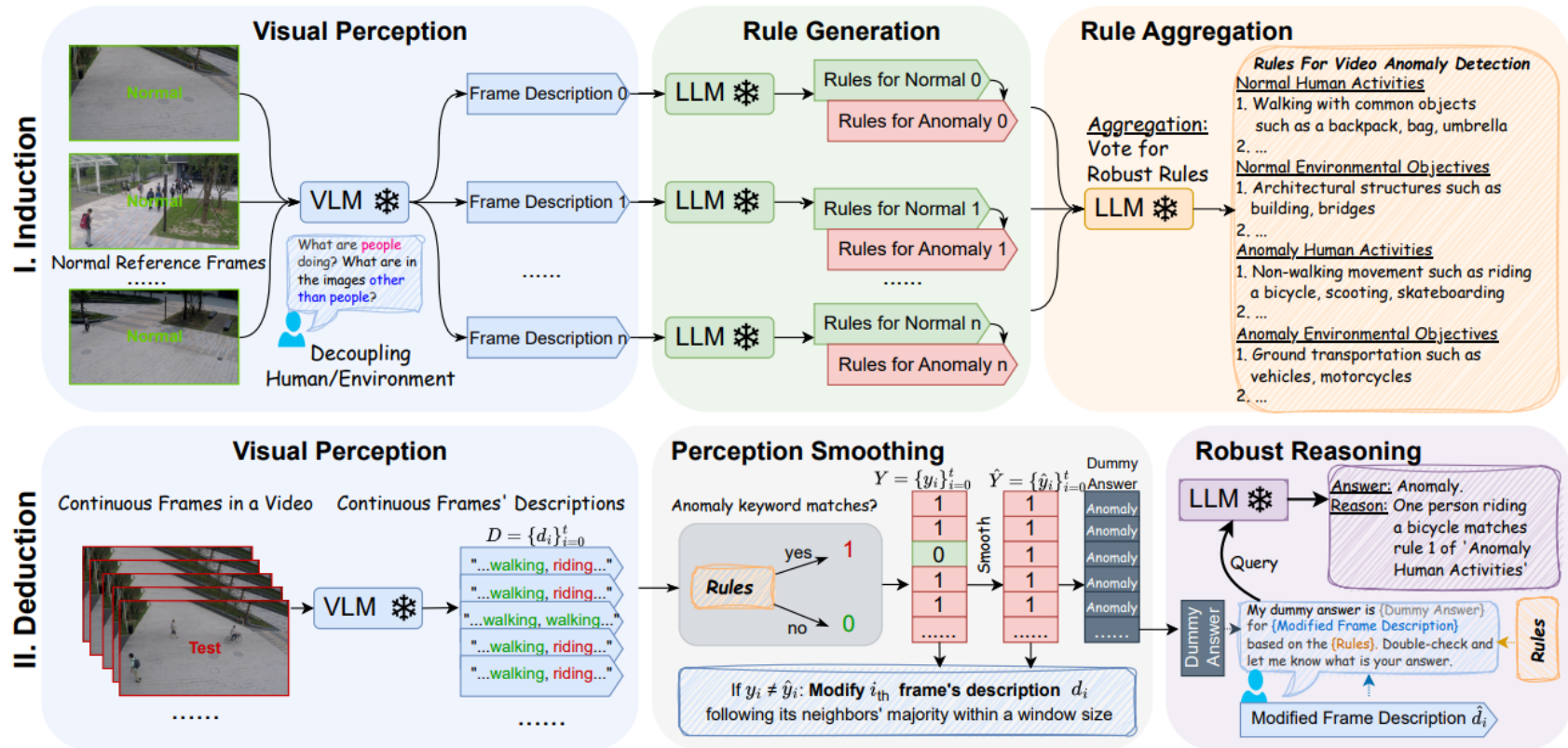- Here GPT-4V mistakenly treats "skateboarding" as normal.

# Method

- **Induction**: Learn rules from few-shot normal reference frames
  - Few-normal shot prompting (no training needed)
- **Deduction**: Detect anomalies based on the rules
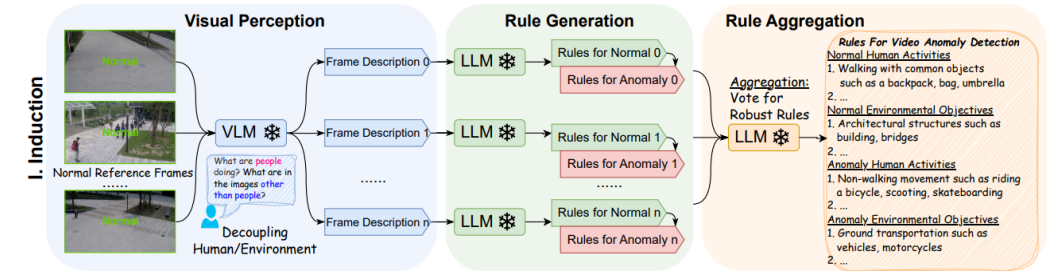  - Correctly identifying "skateboarding" as an anomaly



a) Traditional VAD    b) Ask LLM directly    c) Ask LLM with the rules learned in induction stage (Ours)

# Method

- **Induction (derive rules)**: Use the **few** available normal data as references to derive a set of rules. **Prompting** method without model weight training.

- **Deduction (inference)**: Perform VAD and **explain** detection results according to the induced rules.

# Example (Induction)



A few normal reference frames



A set of rules

```
**Rules for Anomaly Human Activities:**
1. Using any non-walking movement such as riding a bicycle or scooter,
     skateboarding
2. Running or jumping
3. Moving with the usual speed such as very fast or in a hurry
4. Person lying or bending down on the ground
5. Aggressive or unusual movements such as fighting or pushing
6. Loitering around the trash bin


**Rules for Anomaly Environmental Objects:**
1. Missing, removed, or ajar manhole covers
2. Ground transportation such as vehicles, vans, bicycles.
3. Unattended bags or objects (security threat)
4. Unusual objects (vehicles, large unattended luggage, etc.)
5. Unauthorized posters or signs


**Rules for Normal Human Activities:**
1. Walking alone or with another person
2. Standing and using a mobile device
3. Walking on the sidewalk
4. Walking towards or away from a building
5. Entering a building
6. Standing near a trash bin
7. Waiting at traffic lights


**Rules for Normal Environmental Objects:**
1. Manhole covers on the ground or street
2. Trees along the street or Plants present
3. Trash bin in vicinity or placed beside the street
4. Posters on glass door
```
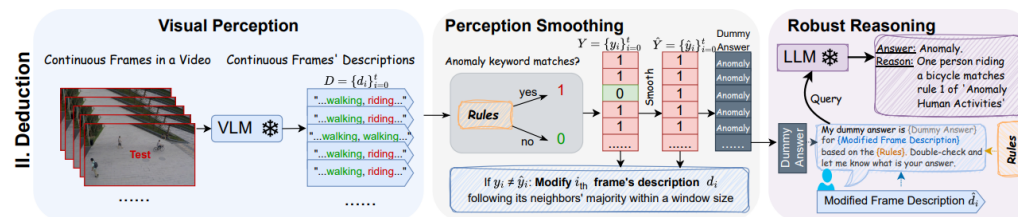
# Example (Deduction)



## Query samples



$\hat{y}_i^* =$ The third person is riding a bicycle, which matches rule number 1 under Anomaly Human Activities.

The bicycle matches rule number 2 under Anomaly Environmental Objects.

Answer: anomaly.

## Output reasoning

$\hat{y}_j^* =$ The human activities present are two people walking. The rules for normal human activities include walking alone or with another person. So, the rule matching is:
1. Walking alone or with another person (normal)

The environmental objects present are two manhole covers on the ground, a trash bin, and some trees and plants. So, the rules matching are:
1. Manhole covers on the ground or street (normal)
2. Trees along the street or Plants present (normal)
3. Trash bin in vicinity or placed beside the street (normal)

Answer: normal.

# Comparison with Conventional VAD Models

- **Induction**: CogVLM-17B & GPT-4

- **Deduction**: CogVLM-17B & Mistral-7B

- **Metric**: AUROC

| Method | Venue | Image Only | Training | Ped2 | Ave | ShT | UB |
|--------|-------|------------|----------|------|-----|-----|-----|
| MNAD [36] | CVPR-20 | ✓ | ✓ | 97.0 | 88.5 | 70.5 | - |
| rGAN [29] | ECCV-20 | ✓ | ✓ | 96.2 | 85.8 | 77.9 | - |
| CDAE [9] | ECCV-20 | ✓ | ✓ | 96.5 | 86.0 | 73.3 | - |
| MPN [30] | CVPR-21 | ✓ | ✓ | 96.9 | 89.5 | 73.8 | - |
| NGOF [50] | CVPR-21 | ✗ | ✓ | 94.2 | 88.4 | 75.3 | - |
| HF2 [25] | ICCV-21 | ✗ | ✓ | **99.2** | 91.1 | 76.2 | - |
| BAF [14] | TPAMI-21 | ✗ | ✓ | 98.7 | 92.3 | 82.7 | 59.3 |
| GCL [56] | CVPR-22 | ✗ | ✓ | - | - | 79.6 | - |
| S3R [53] | ECCV-22 | ✗ | ✓ | - | - | 80.5 | - |
| SSL [49] | ECCV-22 | ✗ | ✓ | 99.0 | **92.2** | 84.3 | - |
| zxVAD [3] | WACV-23 | ✗ | ✓ | 96.9 | - | 71.6 | - |
| HSC [45] | CVPR-23 | ✗ | ✓ | 98.1 | 93.7 | 83.4 | - |
| FPDM [54] | ICCV-23 | ✓ | ✓ | - | 90.1 | 78.6 | 62.7 |
| SLM [43] | ICCV-23 | ✓ | ✓ | 97.6 | 90.9 | 78.8 | - |
| STG-NF [18] | ICCV-23 | ✗ | ✓ | - | - | **85.9** | 71.8 |
| AnomalyRuler-base | - | | ✓ | ✗ | 96.5 | 82.2 | 84.6 | 69.8 |
| AnomalyRuler | - | | ✓ | ✗ | 97.9 | 89.7 | 85.2 | **71.9** |

Two most challenging datasets

# Comparison with LLM-based Approaches

- **Induction**: CogVLM-17B & GPT-4

- **Deduction**: CogVLM-17B & Mistral-7B

- **Metrics**: Accuracy/precision/recall; Doubly-right (RR/RW/WR/WW)

| Method | Accuracy | Precision | Recall |
|---|---|---|---|
| Ask LLM Directly | 52.1 | 97.1 | 6.2 |
| Ask LLM with Elhafsi et al. [12] | 58.4 | **97.9** | 15.2 |
| Ask Video-based LLM Directly | 54.7 | 85.4 | 8.5 |
| AnomalyRuler | **81.8** | 90.2 | **64.3** |

| Method | w. Perception Errors | | | | w/o. Perception Errors | | | |
|---|---|---|---|---|---|---|---|---|
| | RR | RW | WR | WW | RR | RW | WR | WW |
| Ask GPT-4 Directly | 57 | 4 | 15 | 24 | 73 | 3 | 0 | 24 |
| Ask GPT-4 with Elhafsi et al. [12] | 60 | 3 | 15 | 22 | 76 | 2 | 0 | 22 |
| Ask GPT-4V with Cao et al. [8] | 74 | 2 | 7 | 17 | 81 | 2 | 0 | 17 |
| AnomalyRuler | **83** | 1 | 15 | 1 | **99** | 0 | 0 | 1 |

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

Jiacong Xu[1]* Shao-Yuan Lo[2] Bardia Safaei[1] Vishal M. Patel[1] Isht Dwivedi[2]

[1]Johns Hopkins University [2]Honda Research Institute USA

{jxu155, bsafaei1, vpatel36}@jhu.edu {shao-yuan_lo, idwivedi}@honda-ri.com

CVPR 2025

- 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

**Complex reasoning**

## Basic reasoning



**PCB Example: Excessive Solder**

| User | Are there any anomalies 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, specifically the HC-SR04 model. There are no visible anomalies in the image; it appears to be a standard representation 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.

**Macaroni Example: Yellowish Spot**

| User | Do you see any anomalies in the image? |
|---|---|
| GPT-4o | 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. |

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

**In-the-Wild Example: Graffitied Road Sign**

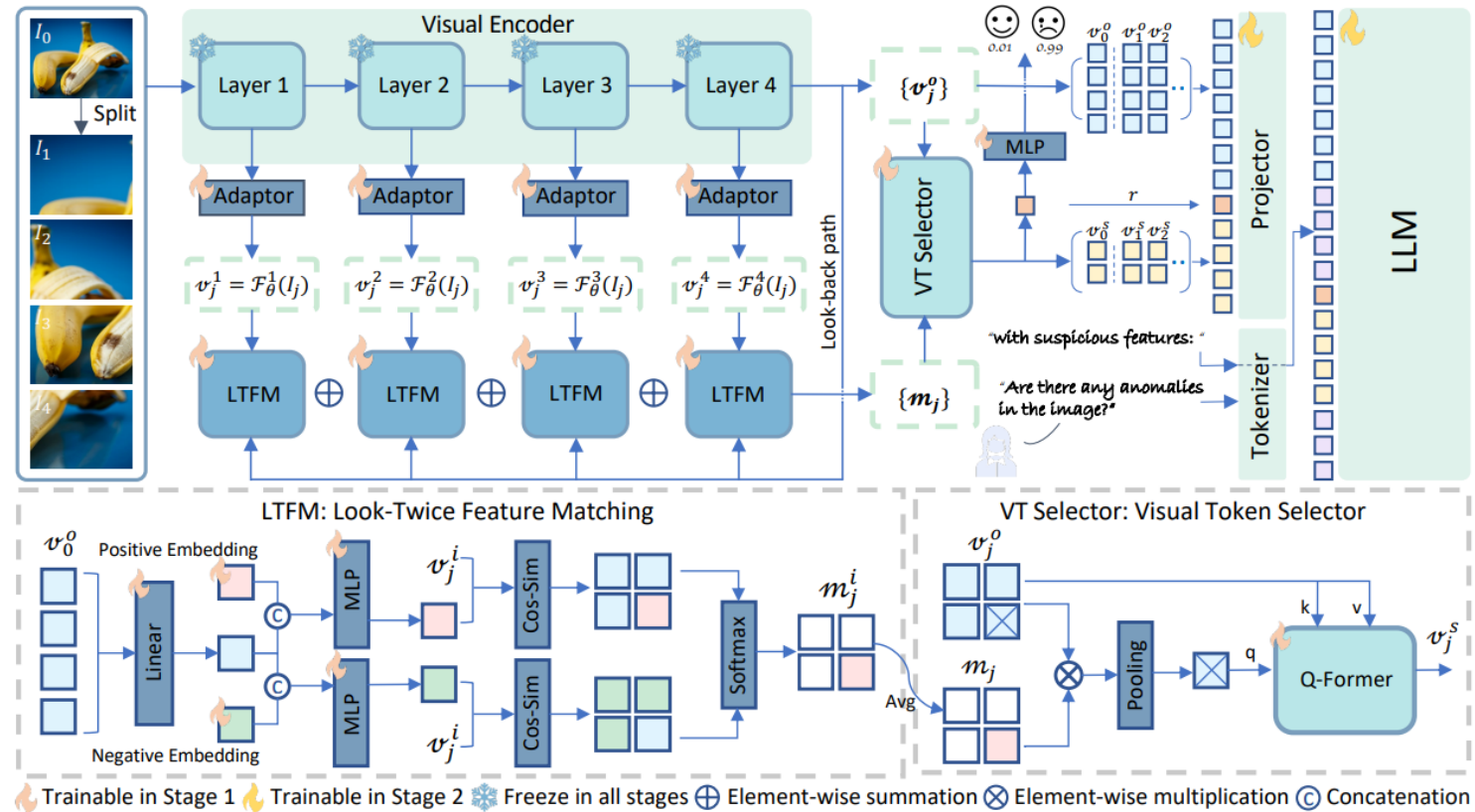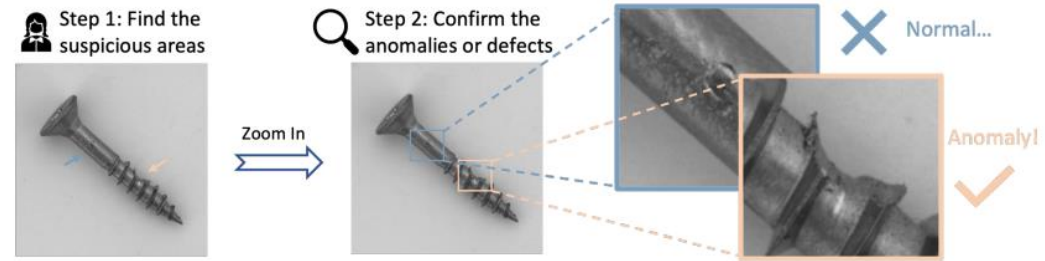| User | Describe the anomaly in the image. |
|---|---|
| Anomaly-OV | The anomaly in the image is the graffiti on the highway sign, which includes black spray paint marks and scribbles over the word 'Seattle' and other parts of the sign. |
| User | What is the reason for the anomaly? |
| Anomaly-OV | The graffiti is likely the result of vandalism, where individuals 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 potentially obstructing the visibility of important information, 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. Additionally, community awareness programs can help reduce vandalism by promoting respect for public property. |

# Method

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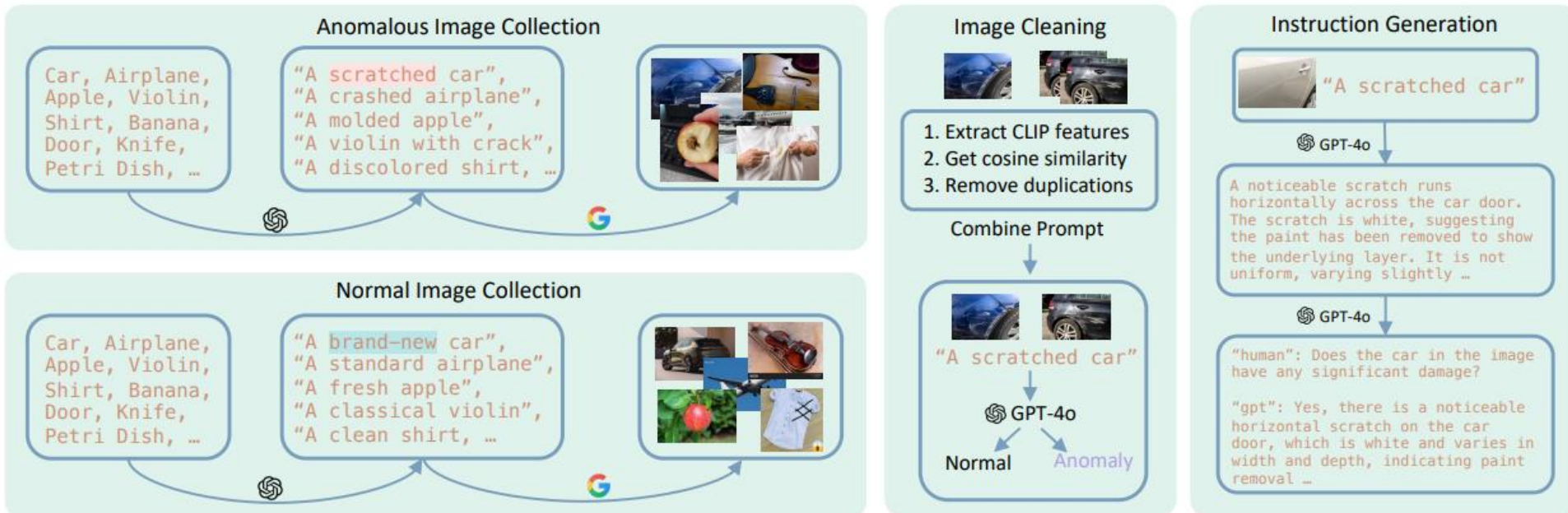- Simulate human's visual anomaly inspection

# 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)

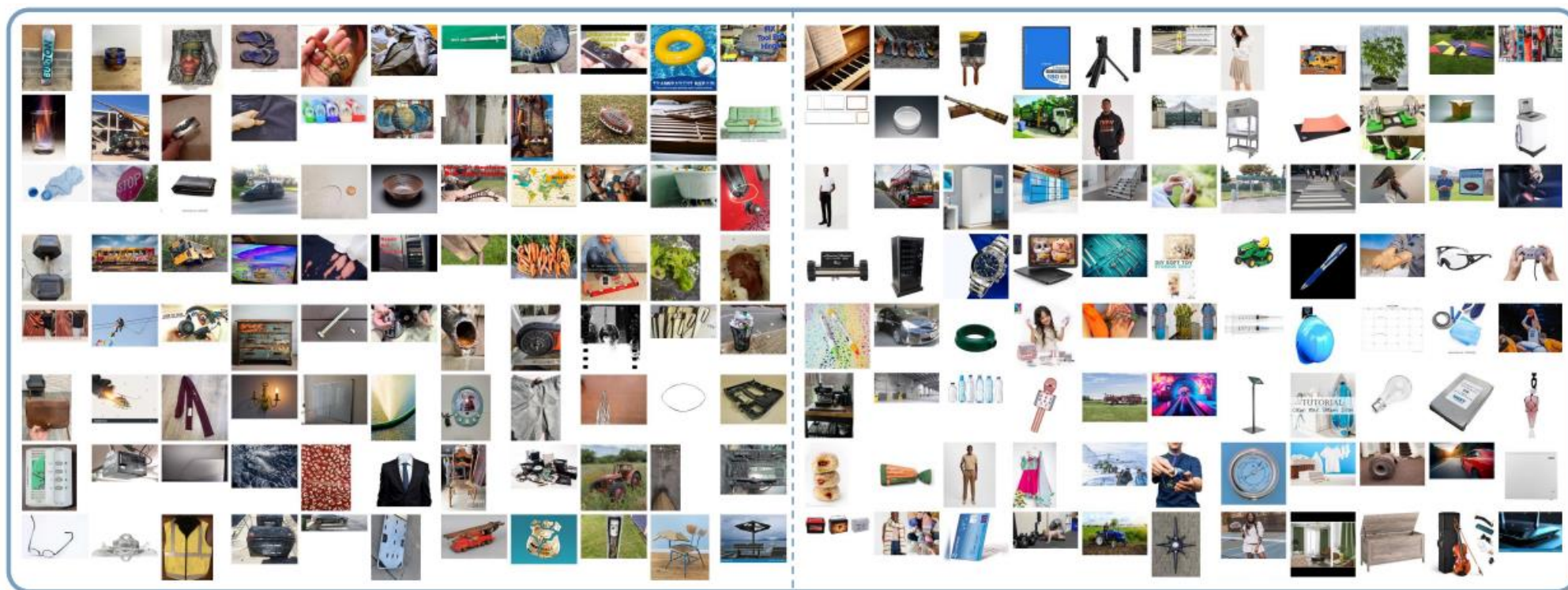- **Training stage 2**: Visual instruction tuning (projector + LLM)

# 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



**WebAD-72k**

[Collected by us via Google Image Search]

MVTec AD

[Bergmann et al. CVPR'19]

BMAD

[Bao et al. CVPR'24]

MVTec-3D AD
Real3D-AD
Anomaly-ShapeNet

[Bergmann et al. 2022] [Liu et al. NeurIPS'23] [Li 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 | |
| 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] | 92.5 | 79.2 | 72.5 | 70.6 | 80.3 | 68.9 | 88.7 | 81.6 | 88.4 | 80.3 |
| AnomalyCLIP [110] | 91.5 | 82.1 | 62.2 | 81.5 | 88.3 | 77.0 | 90.3 | 93.4 | 94.6 | 84.5 |
| AdaCLIP [6] | 89.2 | 85.8 | 64.5 | 79.7 | 88.6 | 76.0 | **94.8** | 91.4 | **97.7** | 85.3 |
| Ours | **94.0** | **91.1** | 72.0 | **83.0** | **89.0** | **81.7** | 93.9 | **97.6** | 95.5 | **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 Reasoning | | | 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-4o [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 | 0.91 | 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 | 0.83 | 5.67 |
| Anomaly-OV-7B | **0.79** | 0.83 | **0.86** | **0.83** | **0.34** | 0.73 | 4.26 | **0.84** | 6.34 |

# Summary

ECCV 2024

Video Anomaly Detection (VAD)

CVPR 2025

Image Anomaly Detection (IAD)

**Conventional**

**MLLMs**

**How to learn**

Full-shot training → Few-shot / Zero-shot

- Prompting pipelines
- Model adaptability

**What to output**

Detection only (anomaly scores) → Detection + Reasoning

- Instruction tuning data
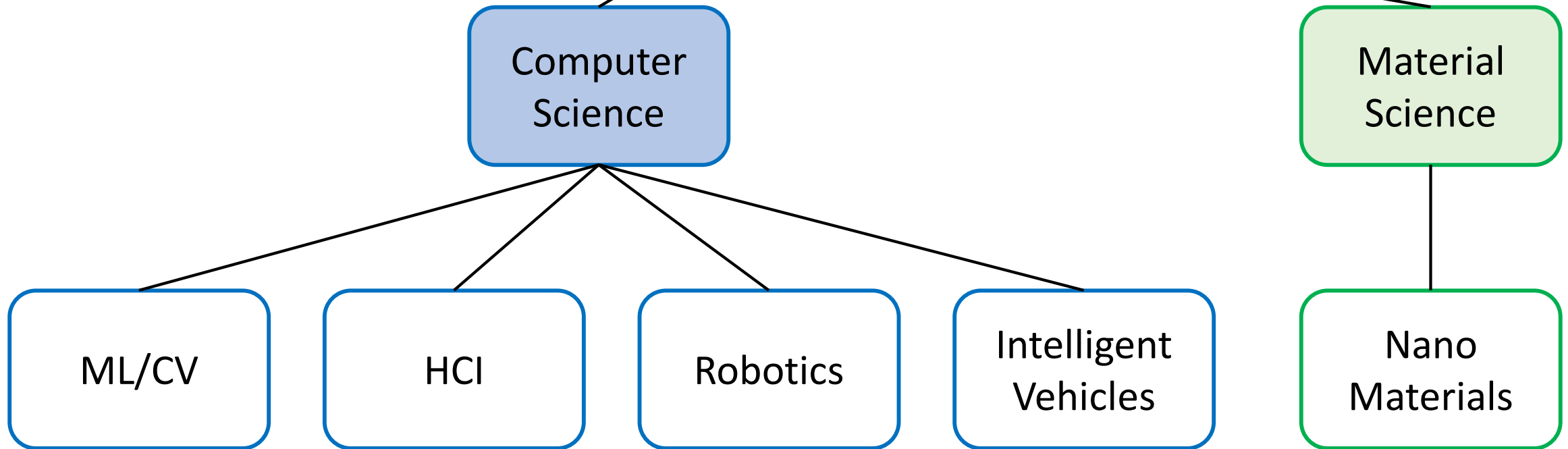- New metrics
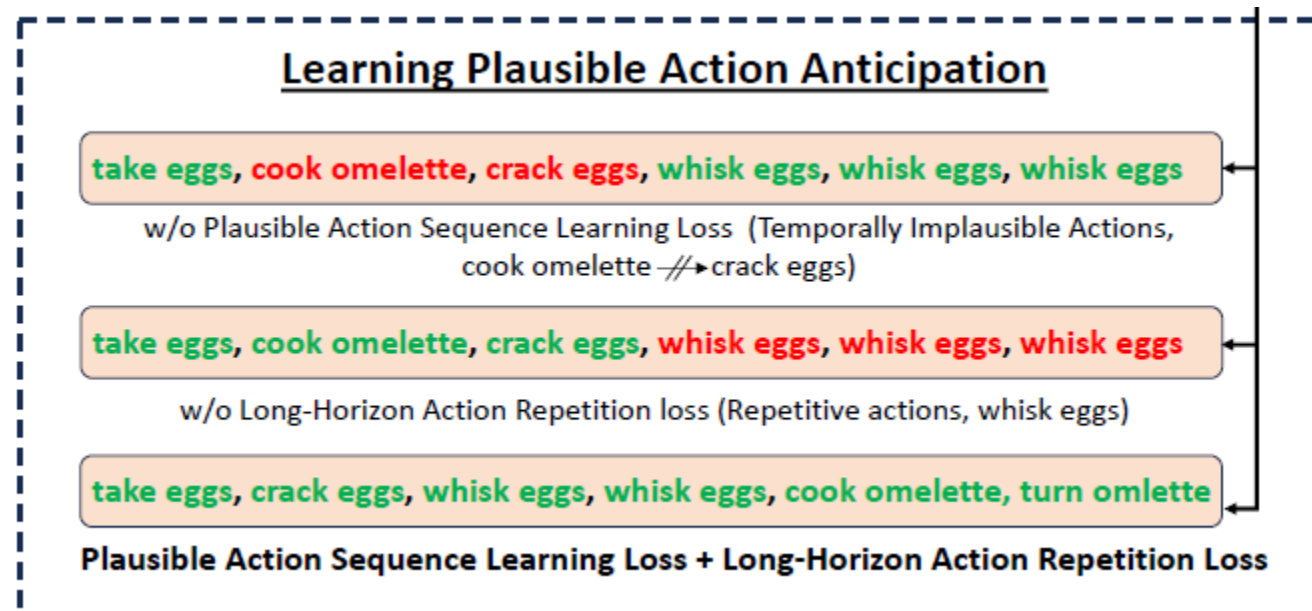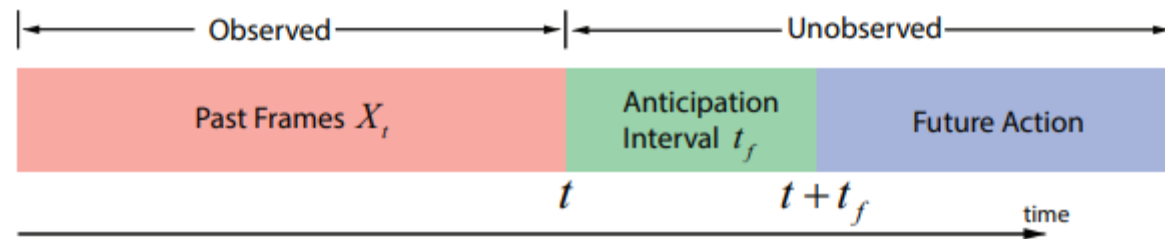- New benchmarks

# Honda Research Institute

# Current Research in HRI ML/CV Team

- **Goal**: Adapt general-purpose MLLMs to domain experts for applications on vehicles or robots.

- MLLM for anomaly detection

- MLLM for action anticipation

- MLLM for affective understanding

- MLLM for Theory-of-Mind

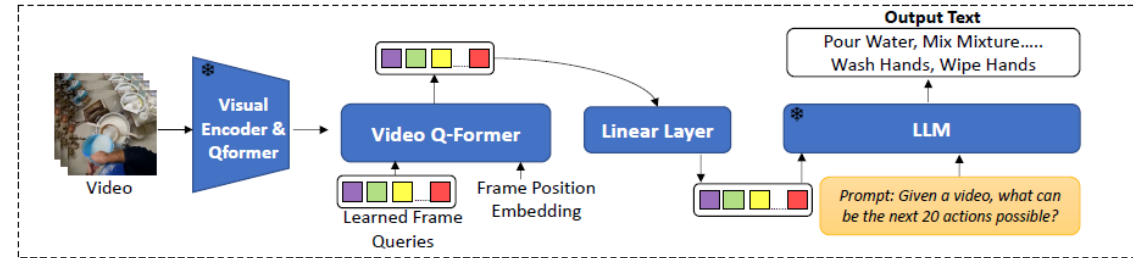- MLLM-based multi-agent collaboration

# MLLM for Action Anticipation

- Action anticipation aims to predict future actions given previous actions.
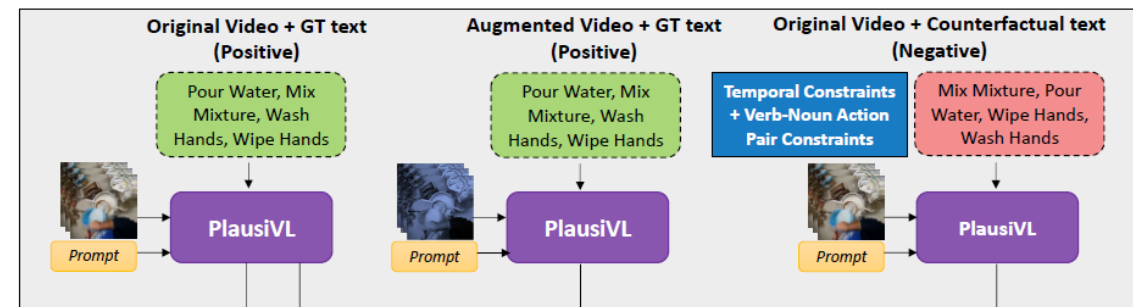


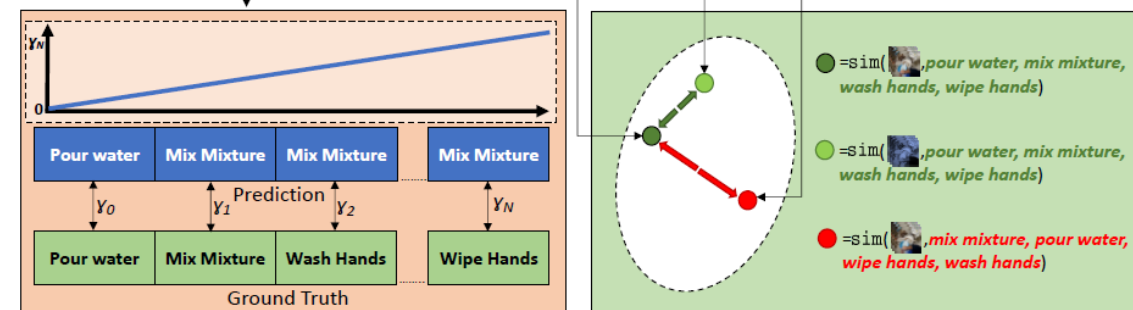CVPR 2024

# MLLM for Action Anticipation

- **Plausible action sequence learning loss**: Help models differentiate between plausible and not plausible action sequences.

- **Long-horizon action repetition loss**: Put a higher penalty on the actions that are more prone to repetition over a longer temporal window.



(a) PlausiVL

(b) Augmentation

(c) Objective Functions and Training

# MLLM for Affective Understanding

- **Task**: Predict and explain viewers' emotional responses to a video.

- Let's start from an example video.
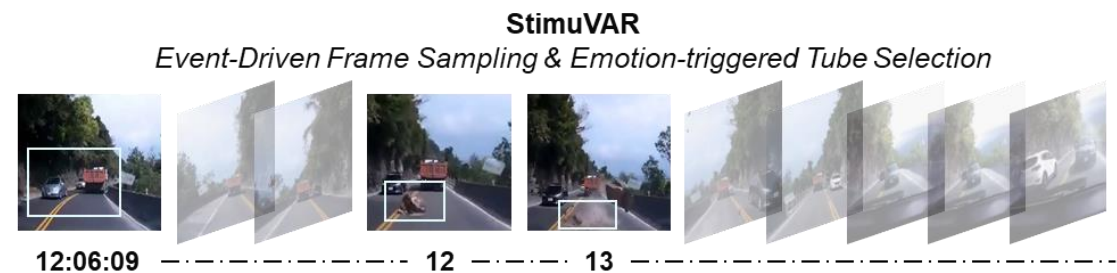
What we see:



What models see:

# MLLM for Affective Understanding

- **Traditional models** and **generic MLLMs** often overlook emotional stimuli.

- **What** are the emotional stimuli?
  - Capturing stimuli could reduce redundancy and improve affective understanding.
  - => Propose **spatiotemporal stimuli-aware mechanisms**

- **Why** this prediction?
  - Interpretability is crucial for model analysis and earning public trust.
  - => Create **VAR visual instruction data** for **affective training**



**Traditional Method**
*Uniform Sampling*

12:06:09 — — 11 — — 14 — — 17 — —

**User:** What is the emotion of the viewers when they watch the video?

**Answer: Boredom** 🐱🐱🐱

**StimuVAR**
*Event-Driven Frame Sampling & Emotion-triggered Tube Selection*

12:06:09 — — 12 — — 13 — —

**User:** What is the emotion of the viewers when they watch the video?

**Answer:** The viewer feels *Surprise* because initially, it is a typical road scene with a red truck and a silver car sets the expectation for a mundane video. However, the sudden appearance of a rock, falling onto the road introduces an unexpected element. This unexpected twist in the video content is likely to evoke surprise in the viewer

*HRI hires research interns!*
*https://usa.honda-ri.com/intern-positions*

About 5 scientists in the ML/CV Team

CVPR 2025          x3

NeurIPS 2024      x2

ECCV 2024          x3

CVPR 2024          x2

*Candidates with at least one first-authored paper at a top conference have a stronger chance.*
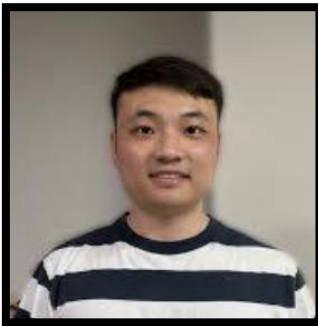
Yuchen Yang

Kwonjoon Lee

Behzad Dariush

Yinzhi Cao

Jiacong Xu

Bardia Safaei

Vishal M. Patel

Isht Dwivedi