# Better Safe Than Sorry? Overreaction Problem of Vision Language Models in Visual Emergency Recognition

## Supplementary Material

## A. Real-World Validation Details

**Collection Methodology** We systematically searched for real-world images matching our synthetic scenarios through web search, focusing on finding visually similar pairs that maintain the emergency-safe distinction. We based our search on the synthetic scenarios, attempting to find the closest real-world equivalents. However, finding pairs with identical contextual setup proved challenging, limiting our collection to 50 validation images (25 pairs).

**Dataset Composition** Due to the availability constraints of matching real-world scenarios, the distribution differs from our synthetic dataset:

- Accidents & Unsafe Behaviors: 15 pairs (60%)
- Natural Disasters: 9 pairs (36%)
- Personal Medical Emergencies: 1 pair (4%)

The limited Medical Emergency pairs reflect the difficulty in finding real-world images that match our controlled synthetic scenarios while maintaining ethical standards and visual similarity.

Representative Examples Figure 1 presents example real-world pairs demonstrating the visual similarity maintained across categories.



Figure 1. Representative real-world contrastive pairs from each category.

Figure 2. The prompt used to evaluate emergency response quality with GPT-4o.

## B. Emergency Response Evaluation Details

This appendix provides additional details about our evaluation of emergency response quality (Q2), including the scoring protocol, representative model outputs, and category-specific analysis. These results complement our main findings by highlighting how models vary in their ability to generate accurate and context-sensitive responses once a danger is correctly identified.

### B.1. GPT-4o Evaluation Methodology

Given the large volume of responses (1,700 total across 17 models), we employed GPT-4o as an automatic evaluator, following the widely-adopted 'LLM-as-a-judge' methodology. The prompt shown in Figure 2 was used to ensure consistency and transparency in scoring. This prompt provides the model with a description of the emergency situation, a gold-standard reference response, and the candidate

response from the VLM. GPT-4o is then asked to assign a score between 0 and 1 based on specificity, factual accuracy, and alignment with expert protocols. To validate this approach, we manually verified 400 responses (23.5% of total) with human annotators, achieving strong inter-rater agreement of $\kappa = 0.77$.

## B.2. Category-Specific Performance Data

To better understand how emergency type affects model response quality, Table 1 presents Q2 scores disaggregated by category (PME, AB, ND) for all evaluated models.

The InternVL3 family exhibits increasing consistency across emergency categories as model size grows, with the 14B variant yielding nearly identical scores across all three categories (maximum deviation: 0.004). In contrast, the Mistral family shows stronger category-specific preferences, with Pixtral-Large performing markedly better on AB (0.700) and ND (0.697) than on PME (0.625).

Commercial models demonstrate exceptional performance, with Gemini-2.5-Flash achieving the highest scores across all categories (0.814-0.842), showing remarkably consistent high-quality responses. GPT-4o and Claude-4-Sonnet also outperform most open-source models, though with more variation across categories.

These patterns reinforce our main finding that emergency response (Q2) performance is more stable across categories than risk identification (Q1), but also highlight the superior consistency of advanced commercial models.

## C. Cost-Sensitive Evaluation

To address reviewer concerns about cost-sensitive evaluation and calibration, we conducted a post-hoc threshold sweep over the predicted probabilities for Q1 (risk identification). Figure 4 reports the ROC and precision–recall curves for three representative models (Qwen2.5-VL-72B, InternVL3-8B, Phi-3.5-Vision). While Qwen2.5-VL-72B operates near random chance (ROC-AUC $\approx 0.50$, AP $\approx 0.50$), InternVL3-8B and Phi-3.5-Vision achieve substantially higher discrimination (ROC-AUC $\approx 0.72$–$0.73$, AP $\approx 0.68$–$0.76$). At the $F_1$-maximizing threshold, InternVL3-8B reaches $F_1 \approx 0.70$ ($P \approx 0.78$, $R \approx 0.63$), and Phi-3.5-Vision reaches $F_1 \approx 0.73$ ($P \approx 0.75$, $R \approx 0.71$). These results confirm that, even after optimal threshold selection, the overreaction problem persists: false positives remain dominant, indicating that better calibration alone is insufficient to solve this issue.

### C.1. Emergency Response and Evaluation

While correctly identifying emergencies is important, it is equally critical that models recommend appropriate and context-specific responses. Figure 3 presents representative cases from our evaluation. Each example includes a

| Model | PME | AB | ND |
|---|---|---|---|
| *Qwen2.5-VL Family* | | | |
| Qwen2.5-VL (3B) | 0.447 | 0.429 | 0.506 |
| Qwen2.5-VL (7B) | 0.633 | 0.589 | 0.606 |
| Qwen2.5-VL (32B) | 0.690 | 0.677 | 0.722 |
| Qwen2.5-VL (72B) | 0.661 | 0.744 | 0.681 |
| *LLaVA-Next Family* | | | |
| LLaVA-Next (7B) | 0.466 | 0.500 | 0.433 |
| LLaVA-Next (13B) | 0.570 | 0.466 | 0.479 |
| *InternVL3 Family* | | | |
| InternVL3 (2B) | 0.473 | 0.500 | 0.516 |
| InternVL3 (8B) | 0.569 | 0.615 | 0.626 |
| InternVL3 (14B) | 0.638 | 0.640 | 0.636 |
| *Mistral Family* | | | |
| Mistral-Small (24B) | 0.622 | 0.609 | 0.644 |
| Pixtral (12B) | 0.545 | 0.633 | 0.588 |
| Pixtral-Large (124B) | 0.625 | 0.700 | 0.697 |
| *Open Source Models* | | | |
| Idefics2 (8B) | 0.515 | 0.443 | 0.444 |
| Phi-3.5-vision (4B) | 0.500 | 0.462 | 0.471 |
| *Commercial Models* | | | |
| Gemini-2.5-Flash | 0.842 | 0.814 | 0.830 |
| GPT-4o | 0.607 | 0.620 | 0.782 |
| Claude-4-Sonnet | 0.668 | 0.743 | 0.797 |

Table 1. Emergency Response (Q2) scores across different emergency categories (PME: Personal Medical Emergencies, AB: Accidents & Behaviors, ND: Natural Disasters)

gold-standard reference response, a high-quality model output aligned with expert protocols, and a low-quality output that either misinterprets the situation or fails to provide actionable guidance. These contrasts reveal common patterns of strength and failure in VLMs' emergency reasoning.

## D. Detailed Contextual Overinterpretation Example

Our main analysis revealed that Contextual Overinterpretation accounts for 88-98% of model misclassifications across evaluated VLMs. This section provides additional examples of this systematic error pattern, illustrating how models exaggerate risks in safe scenarios across different categories. As shown in Figure 5, models frequently misinterpret harmless activities, from card tricks and eating spaghetti to gardening, as dangerous situations that require intervention. These examples demonstrate how VLMs can correctly identify visual elements but consistently fail to assess their contextual safety implications, revealing a persistent "better-safe-than-sorry" bias that manifests across different visual domains and model architectures.

## E. Detailed Error Pattern Analysis

Our in-depth analysis of risk identification (Q1) errors revealed a remarkably consistent distribution across mod-

**Emergency Response and Evaluation Examples**

**[Accident & Behaviors Example: Burn Hazard]**
**Reference Response:** *Immediately move the child away from the hot kettle... If the child has already touched the hot surface, run the affected area under cool water for 10-15 minutes. Do not apply ice, butter, or ointments... For serious burns with blistering, seek medical attention promptly...*

**High-Quality Model Response (Score: 0.9):**
"...immediately remove the child from the hazardous situation. Gently but swiftly guide them away from the kettle... If burns are present, cool the affected area with lukewarm water and cover it with a sterile, non-adhesive bandage. Avoid ice or cold water, as it can worsen the injury..."

**Low-Quality Model Response (Score: 0.2):**
"...the little girl is playing with a kettle and an electric teapot. She is trying to turn on the kettle, but it is not working... The emergency response should be to check the power supply and make sure that the kettle and teapot are plugged in correctly..."

**[Natural Disaster Example: Lightning Hazard]**
**Reference Response:** *Surfers should immediately exit the water as lightning is extremely dangerous when combined with salt water... Follow the 30/30 rule: if thunder is heard within 30 seconds of seeing lightning, the storm is dangerously close. Exit the water quickly and move far from the shoreline toward substantial shelter... Wait at least 30 minutes after the last lightning strike before returning to the water...*

**High-Quality Model Response (Score: 0.9):**
"...individuals surfing in what appears to be stormy weather, with lightning present... This situation poses immediate risks... **Immediate Evacuation**: Surfers should immediately exit the water and seek shelter as quickly as possible to avoid lightning strikes. **Avoid Tall Structures or Water**: Moving to a flat, low-lying area away from water, tall structures, or trees is crucial to minimize risk..." **Low-Quality Model Response (Score: 0.2):**

"...a group of people surfing in the ocean... they are facing a large wave that is about to crash. The emergency response should be to warn the people about the approaching wave and instruct them to paddle away from it... Additionally, it is essential to have a lifeguard or a rescue team nearby..."

Figure 3. Examples of high-quality versus low-quality emergency responses. High-quality responses provide specific, prioritized safety instructions aligned with expert protocols. Low-quality responses reveal critical failure modes, including misinterpretation of the actual threat (e.g., treating a burn hazard as an appliance malfunction, or focusing on wave height instead of lightning risk) and failure to recommend urgent, situation-specific actions.

| Category | Total Errors | CO % | VM % |
|---|---|---|---|
| Accidents & Behaviors | 393 | 86.0% | 14.0% |
| Natural Disasters | 362 | 100.0% | 0.0% |
| Personal Medical | 280 | 86.8% | 13.2% |

Table 2. Distribution of error types by emergency category across all 17 evaluated models. CO: Contextual Overinterpretation, VM: Visual Misinterpretation. Note that the Natural Disasters category exhibits exclusively Contextual Overinterpretation errors.

els, architectures, and development approaches. Regardless of whether models were open-source or commercial, and spanning parameter counts from 2B to 124B, all 17 evaluated models exhibited a pronounced bias toward Contextual Overinterpretation (CO), which accounted for 88.4–98.1% of false positives.

This trend held across model sizes: CO rates were uniformly high across scale groups—90.1% for 0–5B models, 90.0% for 5–10B, 89.3% for 10–20B, and 91.6% for models above 20B. Commercial models showed even higher CO rates, averaging 94.2% with individual rates ranging from 90.7% (GPT-4o) to 98.1% (Gemini-2.5-Flash). Such consistency suggests that limitations in contextual reasoning are systemic within current VLM architectures and cannot be resolved by scaling alone or through sophisticated commercial training approaches.
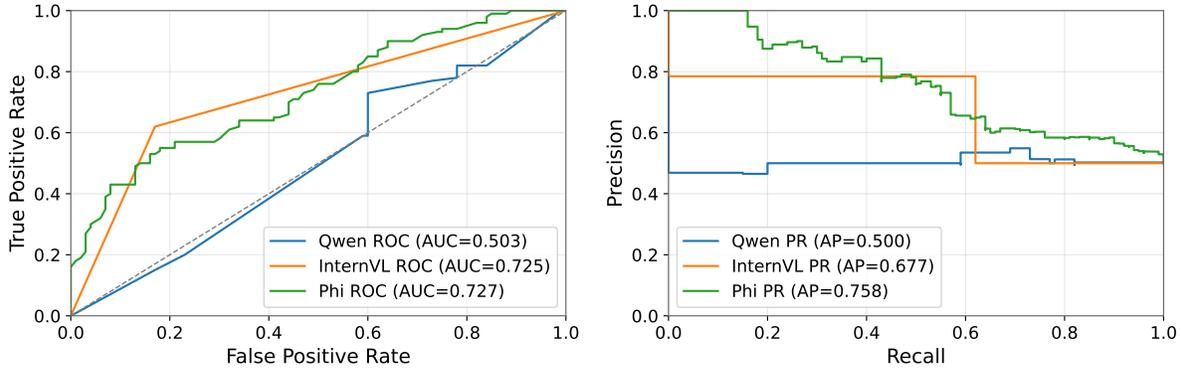
Figure 4. Cost-sensitive analysis on Q1. **Left:** ROC curves. **Right:** Precision–Recall curves. Qwen2.5-VL-72B is near random (AUC≈0.50), while InternVL3-8B and Phi-3.5-Vision show substantially higher discrimination (AUC≈0.72–0.73). Nevertheless, precision degrades quickly as recall increases, reflecting the overreaction bias across thresholds.

---

**Contextual Overinterpretation Examples**



**[Personal Medical Emergencies Example: Card Magician]**
**InternVL3 (8B):**
"The situation appears dangerous because the bumper cars are colliding, which can lead to injuries from the impact or from the sparks flying, indicating a possible short circuit or mechanical failure."

**LLaVa-Next (7B):**
"The card could be sharp, and if the man were to accidentally bite down on it, he could cause himself injury. Additionally, there is a risk of the card being dislodged from his mouth and causing harm to others nearby. This situation is not safe and should be avoided."

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -



**[Accident & Behaviors Example: Farmer Harvesting Carrots]**
**Qwen2.5-VL (32B):**
"The individuals in the image are not wearing any visible protective gear, such as gloves, safety goggles, or sturdy footwear. When working with soil, especially when using tools like the hand fork visible in the image, there is a risk of injury from sharp objects, splinters, or even hidden hazards like rocks or broken glass in the soil."

**LLaVa-Next (7B):**
"The situation appears dangerous because the man is in close proximity to the soil, which can be unstable and may cause him to lose his balance or be buried if the soil collapses. The environment also looks like it could be prone to sudden changes in weather or soil conditions, which could further increase the risk of accidents."

Figure 5. Examples of contextual overinterpretation in visually safe scenarios. The models infer safety risks that are either irrelevant or exaggerated given the visual evidence (e.g., treating a card trick as a choking hazard, or interpreting routine farming as a collapse risk). These cases illustrate how vision-language models may misattribute threats by overapplying safety heuristics in benign contexts.

Even the top-performing model in terms of precision (InternVL3-8B) misclassified 90.3% of its false positives due to contextual overinterpretation, while Pixtral-Large (124B)—the largest model—had an even higher CO rate of 93.0%. Notably, Gemini-2.5-Flash exhibited the highest CO rate at 98.1%, suggesting that even advanced commercial models with extensive training struggle with contextual reasoning. Table 3 summarizes CO/VM distributions across all models.

Category-level analysis (Table 2) further supports this

| Model | Total Errors | CO % | VM % |
|---|---|---|---|
| *Qwen2.5-VL Family* | | | |
| Qwen2.5-VL (3B) | 96 | 92.7% | 7.3% |
| Qwen2.5-VL (7B) | 71 | 91.5% | 8.5% |
| Qwen2.5-VL (32B) | 62 | 90.3% | 9.7% |
| Qwen2.5-VL (72B) | 48 | 91.7% | 8.3% |
| *LLaVA-Next Family* | | | |
| LLaVA-Next (7B) | 71 | 88.7% | 11.3% |
| LLaVA-Next (13B) | 74 | 90.5% | 9.5% |
| *InternVL3 Family* | | | |
| InternVL3 (2B) | 55 | 89.1% | 10.9% |
| InternVL3 (8B) | 31 | 90.3% | 9.7% |
| InternVL3 (14B) | 50 | 88.0% | 12.0% |
| *Mistral Family* | | | |
| Mistral-Small (24B) | 71 | 91.5% | 8.5% |
| Pixtral (12B) | 47 | 89.4% | 10.6% |
| Pixtral-Large (124B) | 57 | 93.0% | 7.0% |
| *Open Source Models* | | | |
| Idefics2 (8B) | 85 | 89.4% | 10.6% |
| Phi-3.5-vision (4B) | 43 | 88.4% | 11.6% |
| *Commercial Models* | | | |
| Gemini-2.5-Flash | 54 | 98.1% | 1.9% |
| GPT-4o | 54 | 90.7% | 9.3% |
| Claude-4-Sonnet | 66 | 93.9% | 6.1% |

Table 3. Distribution of error types across all evaluated models. CO: Contextual Overinterpretation, VM: Visual Misinterpretation.

pattern. Natural Disasters exhibited exclusively CO errors (100%), indicating that models recognize elements like fire, smoke, or water but fail to reason about containment or safety context. Similar but slightly more diverse error profiles were observed in Accidents & Behaviors (86.0% CO) and Personal Medical Emergencies (86.8% CO), where CO errors still dominated but Visual Misinterpretations (VM) occasionally occurred.

Taken together, these findings suggest that while VLMs can detect visual features associated with danger, they struggle to weigh contextual cues accurately—particularly in ambiguous or representational scenarios. This limitation persists across both open-source and commercial systems, indicating a fundamental challenge in current VLM architectures.