

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

Dasol Choi^{1,2} Seunghyun Lee¹ Youngsook Song^{3*}

¹Yonsei University ²MODULABS ³Lablup Inc.

{dsaolchoi, lutris}@yonsei.ac.kr yssong@lablup.com

Abstract

Vision-Language Models (VLMs) have shown capabilities in interpreting visual content, but their reliability in safety-critical everyday life scenarios remains insufficiently explored. We introduce VERI (Visual Emergency Recognition Dataset), a diagnostic benchmark comprising 200 images organized into 100 contrastive pairs. Each emergency scene is paired with a visually similar but safe counterpart through human verification and refinement. Using a two-stage evaluation protocol—risk identification and emergency response—we assess 14 VLMs (2B to 124B parameters) across medical emergencies, accidents, and natural disasters. Our analysis reveals an “overreaction problem,” where models accurately identify genuine emergencies (70-100% success rate) but produce high false-positive rates, misclassifying 31-96% of safe situations as dangerous. Ten safe scenarios were universally misclassified by all models regardless of scale. This “better-safe-than-sorry” bias primarily results from contextual overinterpretation (88-93% of errors), challenging VLM reliability in safety-critical applications. These findings highlight fundamental limitations in current VLM architectures, which persist despite increased model scale. Our results demonstrate an urgent need for strategies specifically improving contextual reasoning in ambiguous visual situations. The consistently low performance of the model indicates that these data serve effectively as a diagnostic data set.

Content Warning: This paper contains images and descriptions of emergency situations.

1. Introduction

Vision-Language Models (VLMs) have made remarkable progress in understanding visual content, advancing from

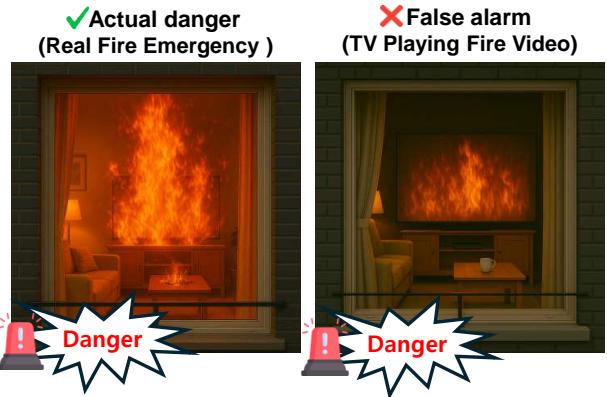


Figure 1. The overreaction problem in VLMs: correctly identifying actual emergencies (left) while misclassifying visually similar safe scenarios as dangerous (right), similar to human misperceptions of TV fire videos as real fires in 2023.

simple object recognition to sophisticated scene understanding and contextual reasoning [4, 10, 17]. These models now power a wide range of applications from content moderation to assistive technologies [26]. However, their reliability in safety-critical scenarios remains underexplored. This gap raises a critical question with real-world implications: To what extent can current VLMs distinguish between genuine emergencies and visually similar but safe scenarios? To investigate this question systematically, we utilized generated synthetic (fake) data, which allows us to create controlled comparisons between dangerous and benign situations while maintaining visual similarity.

False visual perceptions of emergencies can lead to costly resource mobilization. In October 2023, firefighters were dispatched to two separate incidents in New York and Seoul. In both cases, they discovered that high-definition fireplace videos playing on television screens had been mistaken for actual fires [15, 23]. Such false alarms carry substantial costs. Each incident costs approximately \$1,000-2,500, with annual costs exceeding \$1 billion in the United States alone [19]. As VLMs power applications from

*Corresponding Author

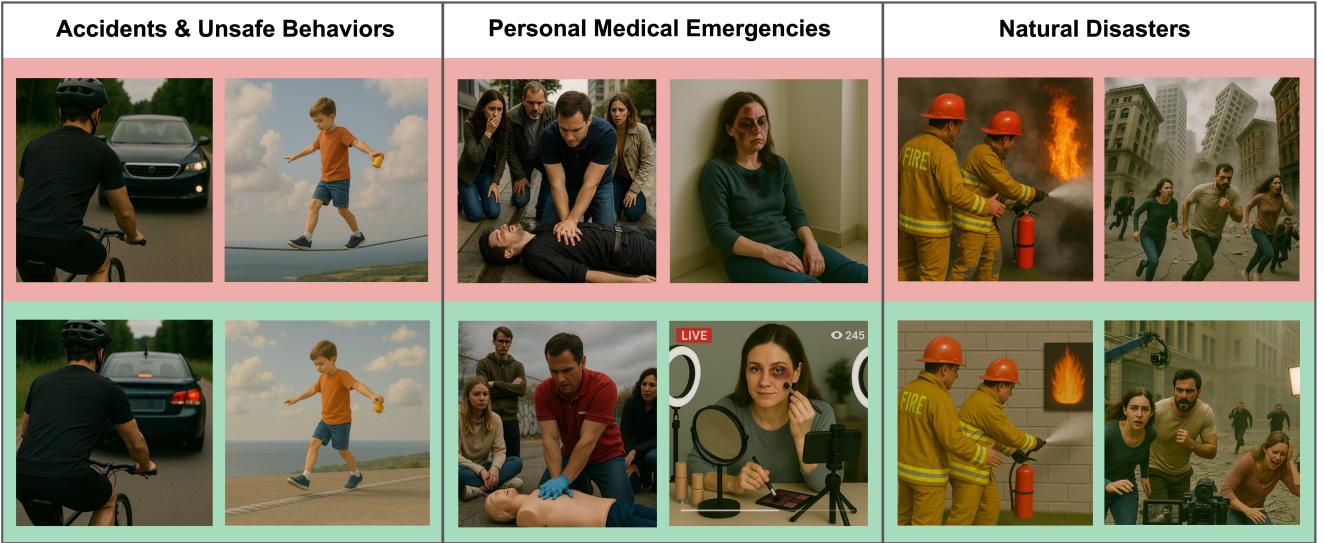


Figure 2. Examples from the VERI dataset showing contrastive pairs across three categories. Top row (red background): genuine emergency situations. Bottom row (green background): visually similar but safe scenarios. Each pair maintains visual similarity while representing different semantic meanings—one requiring intervention and the other representing safe activities. These pairs enable evaluation of VLMs’ ability to make safety distinctions despite visual similarities.

smart home monitoring to CCTV surveillance, understanding their tendency to “overreact” becomes crucial [24]. As VLMs increasingly power applications from smart home monitoring to CCTV surveillance, understanding their tendency to ‘overreact’ becomes crucial [24]. Current research on safety-critical VLMs has primarily focused on specialized domains. These include medical diagnosis systems [6], autonomous driving perception [7], and industrial robotics safety mechanisms [8]. Our work addresses a distinct gap: everyday emergency recognition across domains. This task requires different contextual reasoning capabilities—distinguishing between visually similar scenarios based on subtle contextual cues rather than domain-specific features. Existing benchmarks typically evaluate models on isolated images, rarely using contrastive pairs that test the ability to distinguish between visually similar but semantically distinct situations.

To address these gaps, we introduce VERI (Visual Emergency Recognition Dataset), a diagnostic benchmark of 200 images comprising 100 contrastive pairs. Inspired by compact yet influential diagnostic benchmarks such as Winogender (720 examples) [18], WSC (60 examples) [13], and PAIRS (200 images) [9], VERI prioritizes diagnostic value over quantity. We employ controlled image pairs to evaluate VLMs’ ability to distinguish between emergency and non-emergency situations. Each pair underwent over 5 rounds of iterative refinement and verification to balance visual similarity with semantic distinction.

Our work makes the following contributions:

- We introduce VERI, a benchmark with a two-stage proto-

col assessing risk identification and action suggestion in emergency scenarios.

- Through evaluation of 14 VLMs (2B-124B parameters), we reveal a systematic “overreaction problem” persisting across architectures and scales, with 10 safe scenarios misclassified by all models.
- We identify two error patterns (contextual and visual misinterpretations), with contextual errors dominating (88-93%) and varying by emergency category.
- We show that the overreaction problem persists despite model scaling, requiring targeted architectural improvements for contextual emergency assessment.

2. VERI: Visual Emergency Recognition Dataset

2.1. Dataset Design and Taxonomy

Effective emergency detection requires distinguishing between genuine threats and visually similar but safe situations. To address this challenge, the VERI dataset (200 images) has been specifically designed to study risks in everyday life, rather than focusing on specialized domains such as medical diagnosis or industrial safety. This focus on everyday risks is particularly relevant given the widespread adoption of autonomous driving and AI-driven robotics in our daily lives. By establishing this diagnostic benchmark, our study deliberately narrows its research scope to typical daily-life scenarios, where accurate threat detection is becoming increasingly critical. The core design principle of VERI is the contrastive pair approach—each entry consists

Category	Accidents & Unsafe Behaviors	Personal Medical Emergencies	Natural Disasters
Scope	Immediate physical dangers from environment or human action	Urgent health risks to individuals	Large-scale threats affecting multiple people
Example scenarios	Traffic accidents, falls from heights, drowning risks, physical altercations, unsafe tool use	Cardiac arrest, choking, unconsciousness, severe injuries, allergic reactions, seizures	Fires, floods, earthquakes, building collapses, hurricanes, avalanches

Table 1. Taxonomy and examples of emergency situations in the VERI dataset.

of two images with high visual similarity but fundamentally different semantic implications: one depicting a genuine emergency requiring intervention and the other showing a visually similar but safe scenario. As illustrated in Figure 2, these pairs include contrasts such as actual medical emergencies versus training simulations, genuine accidents versus staged scenarios, and real disasters versus controlled environments. This contrastive structure enables direct assessment of a model’s ability to distinguish between superficially similar scenes with different safety implications. We organized VERI into three major categories of emergency situations (Accidents & Unsafe Behaviors, Personal Medical Emergencies, and Natural Disasters), as shown in Table 1. Each category presents unique visual challenges, from subtle physiological distress cues to complex environmental interpretation, as demonstrated in Figure 2.

2.2. Image Creation Process

Creating visually similar but semantically distinct image pairs required careful design to balance realism with controlled visual properties. We employed a multi-stage approach:

Stage 1: Scenario Definition. For each category, we defined specific scenario pairs (e.g., “person requiring CPR” vs. “CPR training on mannequin”), identifying visual elements to preserve or alter for semantic distinction.

Stage 2: Image Generation. Two researchers created prompts for image generation using GPT-4o. We preserved key visual elements while altering critical semantic cues (e.g., real person vs. mannequin, distress vs. educational context).

Stage 3: Refinement and Validation. We iteratively refined images using inpainting until they matched our requirements (averaging 5+ cycles per pair). Three evaluators (two researchers and an independent annotator) validated each image against our criteria. The process continued until *unanimous agreement* was reached on all criteria (Cohen’s $\kappa = 1.00$, $n = 200$ pairs, 600 decisions), eliminating annotation noise. These criteria ensured emergency scenes clearly required intervention, safe scenes were unambiguously non-threatening, and pairs maintained appropriate visual similarity.

Stage 4: Caption Generation. We created descriptive captions for each image pair using GPT-4o, which were then

refined to ensure accuracy.

Our process prioritized diagnostic value over quantity, focusing on expert-guided refinement to effectively challenge VLMs’ contextual reasoning. Because data for everyday emergency recognition did not previously exist, generating a large volume of realistic images by abstracting from existing raw data was impractical. Moreover, for diagnostic benchmarks with specific evaluation goals, the qualitative value and the control of variables are significantly more important than sheer quantity. This principle aligns with other influential diagnostic benchmarks in the literature. For instance, the PAIRS benchmark comprises 200 images, and the study ‘Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena’ (Zheng et al., 2023) evaluated models using only 80 carefully constructed questions across eight fields.

2.3. Query Development

Our evaluation protocol requires two types of queries to assess different aspects of model performance:

Q1: Risk Identification. The first query—“Does the situation in this image appear dangerous or safe?”—is a straightforward binary classification task designed to evaluate models’ ability to distinguish between dangerous and safe situations.

Q2: Emergency Response. For emergency images, we used the query—“Based on this situation, what emergency response should be taken?”—to assess models’ ability to recommend appropriate actions. We created a standardized response bank through a three-step process. First, we used image captions as input to Perplexity to generate initial responses based on authoritative sources. Next, three evaluators verified these responses against original sources. Finally, we compiled critical actions for each emergency type as evaluation criteria.

2.4. Dataset Statistics

The final VERI dataset consists of 100 image pairs (200 total images) distributed across three categories: Accidents & Unsafe Behaviors (35 pairs), Personal Medical Emergencies (33 pairs), and Natural Disasters (32 pairs), as summarized in Table 3. For evaluation, we created 200 binary classification questions (Q1) covering all images and 100 open-ended response questions (Q2) for emergency images only. Each image is accompanied by a detailed caption describ-

Model	Q1: Risk Identification			Q2: Emergency Response	
	Precision	Recall	F1	Score	# Images
<i>Qwen2.5-VL Family</i>					
Qwen2.5-VL (3B)	0.510	1.000	0.676	0.460	98
Qwen2.5-VL (7B)	0.554	0.880	0.680	0.618	88
Qwen2.5-VL (32B)	0.589	0.890	0.709	0.6972	88
Qwen2.5-VL (72B)	0.652	0.900	0.756	0.7000	89
<i>LLaVA-Next Family</i>					
LLaVA-Next (7B)	0.577	0.970	0.724	0.466	95
LLaVA-Next (13B)	0.575	1.000	0.730	0.502	98
<i>InternVL3 Family</i>					
InternVL3 (2B)	0.633	0.950	0.760	0.497	93
InternVL3 (8B)	0.721	0.800	0.758	0.610	80
InternVL3 (14B)	0.658	0.960	0.781	0.638	94
<i>Mistral Family</i>					
Mistral-Small (24B)	0.572	0.950	0.714	0.625	93
Pixtral (12B)	0.654	0.890	0.754	0.594	89
Pixtral-Large (124B)	0.632	0.980	0.769	0.677	96
<i>Other Models</i>					
Idefics2 (8B)	0.528	0.950	0.679	0.463	93
Phi-3.5-vision (4B)	0.620	0.700	0.657	0.471	70

Table 2. Performance evaluation across risk identification (Q1) and emergency response (Q2) tasks. Q1 metrics show models’ ability to distinguish between dangerous and safe situations. Q2 scores reflect the quality of suggested actions for correctly identified emergencies, with # Images indicating the number of emergency images for which the model provided recommendations.

Statistic	Count
Total image pairs	100
Total images	200
Accidents & Unsafe Behaviors pairs	35
Personal Medical Emergencies pairs	33
Natural Disasters pairs	32
Risk identification QA (Q1)	200
Emergency Response QA (Q2)	100
Detailed image captions	200

Table 3. VERI dataset statistics

ing the scene context, key elements, and situational details, providing additional textual information that can be used for multimodal training or analysis. Figure 2 illustrates representative examples from each category, demonstrating both the visual similarity within pairs and the semantic distinction between emergency and safe scenarios.

3. Experimental Settings

3.1. Models

We evaluated 14 open-source VLMs (2B-124B parameters) across different architectures: Qwen2.5-VL [5] (3B, 7B, 32B, 72B), transformer-based models optimized for visual

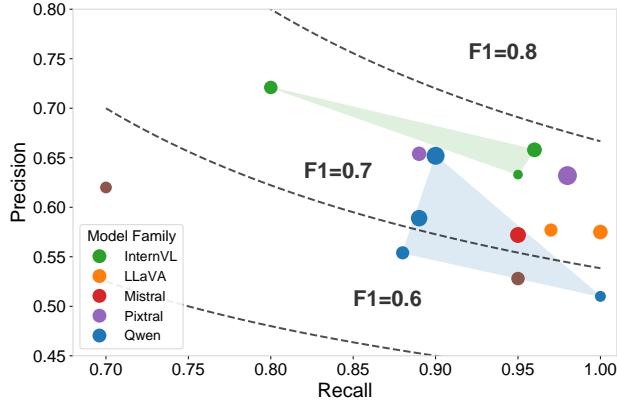
recognition; LLaVA-Next [14] (7B, 13B), integrating CLIP encoders with large language models; InternVL3 [27] (2B, 8B, 14B), using a “ViT-MLP-LLM” architecture; the Mistral family , including Mistral-Small (24B) [3] and Pixtral variants (12B, 124B) [2]; and other architectures represented by Idefics2 (8B) [12] and Phi-3.5-vision (4B) [1]. All models were evaluated using publicly available checkpoints without task-specific fine-tuning to assess zero-shot capabilities.

3.2. Evaluation Protocol

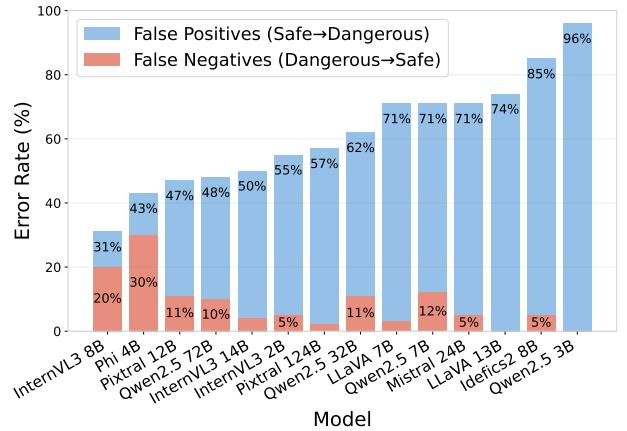
We evaluate model responses using a two-stage protocol:

Q1 (Binary Classification). The model is asked whether the situation in the image is dangerous or safe. It selects one of two choices (A. Dangerous / B. Safe) and provides a brief reasoning. We compute *precision*, *recall*, and *F1 score* based on human-annotated binary ground truth labels.

Q2 (Open-ended Response). For images classified as dangerous by the model and correctly aligned with the ground truth, we further ask what an appropriate emergency response should be. We evaluated these responses using GPT-4o as a judge, providing it with the image caption and our curated gold-standard answers as references. The judge was instructed to score each response on a scale from 0 to



(a) Precision-Recall tradeoff



(b) False Positive vs False Negative rates

Figure 3. Performance analysis of VLMs on emergency detection tasks. (a) Shows the pattern of high recall but lower precision across models, with point size indicating model parameter count. (b) Reveals a consistent “better-safe-than-sorry” bias where safe images are misclassified more frequently than emergencies are missed.

1 based on its alignment with the reference materials and appropriateness for the emergency situation.

4. Results and Analysis

4.1. The Overreaction Phenomenon

To evaluate VLMs’ emergency recognition capabilities, we assessed their performance on both risk identification (Q1) and emergency response (Q2) tasks (Table 2). Our evaluation reveals a consistent pattern: models achieve high recall (0.70-1.00) in identifying dangerous situations, but precision is notably lower (0.51-0.72), indicating a systematic “better-safe-than-sorry” bias—an “overreaction problem.” This pattern persists across model families and parameter scales (2B-124B). Model size does not consistently correlate with improved precision; InternVL3 (8B) achieves the highest precision (0.72), outperforming larger variants. False positive rates (safe images misclassified as dangerous) ranged from 31% to 96%—substantially higher than false negative rates (missed emergencies), which ranged from 0% to 30%. Figure 3(a) illustrates this bias through the precision-recall tradeoff, with models clustering in the high-recall but lower-precision region.

Particularly concerning is that 10% of safe scenarios were misclassified by all 14 evaluated models, typically containing visual elements strongly associated with danger despite clear contextual safety cues. As shown in Figure 3(b), false positives consistently outnumber false negatives across all models. Increasing model size did not reliably mitigate this problem—within the Qwen2.5-VL family, precision improves modestly from 3B (0.51) to 72B (0.65), but recall fluctuates non-linearly. This persistent pattern across architectures and scales suggests the overreaction problem may be embedded in foundational visual un-

derstanding rather than reasoning capacity.

4.2. Emergency Response Evaluation

Beyond identifying emergencies, we evaluated how effectively models suggest appropriate actions for emergency situations. Models demonstrated moderate effectiveness (scores 0.46-0.70) in generating appropriate emergency responses for correctly identified dangerous scenarios.

Unlike risk identification, emergency response quality shows a clearer correlation with model size. We observe consistent scaling benefits: Qwen2.5-VL (0.46→0.70), InternVL3 (0.50→0.64), and Mistral (0.59→0.68). This linear relationship presents an interesting contrast to the overreaction problem: while contextual reasoning for risk assessment doesn’t consistently improve with scaling, procedural knowledge for emergency responses does. This suggests that protocols may be more amenable to parameter scaling than the nuanced contextual judgments required for accurate risk identification.

We found a disconnect between risk identification and response capabilities. InternVL3 (8B) achieved the highest precision in risk identification but only moderate response quality, while models with lower precision often produced better emergency responses. This suggests VLMs develop these capabilities through different mechanisms—procedural knowledge scaling more predictably than contextual reasoning. Even when models correctly identify emergencies, their ability to recommend appropriate interventions remains limited. Examples can be found in the supplementary material.

4.3. Category-Specific Analysis

Our analysis reveals measurable performance variations across emergency categories (Table 4). In risk identifica-



Figure 4. Error patterns in danger assessment: Visual Misinterpretation (blue) vs. Contextual Overinterpretation (red). Examples show how VLMs incorrectly classify safe situations as dangerous through either misperceiving visual elements or exaggerating potential risks in properly perceived scenarios.

tion (Q1), models generally performed best on Natural Disasters (ND, avg F1=0.75) compared to Accidents & Behaviors (AB, 0.73) and Personal Medical Emergencies (PME, 0.68). The PME category shows the highest variance, from exceptional (0.79 for Pixtral-Large) to extremely poor (0.41 for Phi-3.5-vision).

Model scaling effects vary by category. For PME, larger models generally show improved performance (e.g., Qwen2.5-VL improves from 0.68 to 0.73, Pixtral from 0.67 to 0.79), suggesting increased parameters help with subtle physiological cues. For AB and ND categories, we observe a non-linear relationship with size, where mid-sized models sometimes outperform larger counterparts (e.g., InternVL3-8B achieves 0.83 on AB, outperforming the 14B variant). The overreaction problem is most pronounced in the PME category (38% false positive rate) where models frequently misclassified training scenarios as genuine emergencies. Natural Disasters show the lowest false positive rate (27%), likely because environmental containment features provide more distinctive safety indicators.

These findings suggest VLMs struggle most with scenarios where emergency status depends on fine-grained contextual details rather than obvious visual elements. For emergency response quality, Q1 showed a clearer category preference (ND > AB > PME with differences up to 0.07), while Q2 exhibited smaller average differences between categories (approximately 0.03). This suggests that once a model correctly identifies an emergency, its ability to generate appropriate responses is more consistent across different emergency types. For detailed model-specific performance breakdowns across emergency categories, see the supplementary material.

Model	PME	AB	ND
<i>Qwen2.5-VL Family</i>			
Qwen2.5-VL (3B)	0.681	0.673	0.674
Qwen2.5-VL (7B)	0.583	0.737	0.696
Qwen2.5-VL (32B)	0.712	0.697	0.719
Qwen2.5-VL (72B)	0.727	0.764	0.771
<i>LLaVA-Next Family</i>			
LLaVA-Next (7B)	0.729	0.695	0.750
LLaVA-Next (13B)	0.727	0.729	0.733
<i>InternVL3 Family</i>			
InternVL3 (2B)	0.753	0.747	0.781
InternVL3 (8B)	0.593	0.825	0.805
InternVL3 (14B)	0.767	0.761	0.815
<i>Mistral Family</i>			
Mistral-Small (24B)	0.700	0.707	0.736
Pixtral (12B)	0.667	0.767	0.815
Pixtral-Large (124B)	0.790	0.745	0.777
<i>Other Models</i>			
Idefics2 (8B)	0.683	0.693	0.660
Phi-3.5-vision (4B)	0.408	0.707	0.756

Table 4. F1 scores across different emergency categories (PME: Personal Medical Emergencies, AB: Accidents & Behaviors, ND: Natural Disasters)

4.4. Error Patterns and Analysis

Our analysis reveals two primary patterns in the overreaction problem (false positives), where models misclassify safe situations as dangerous, as illustrated in Figure 4:

Visual Misinterpretation Models incorrectly perceive visual elements, failing to distinguish between safe and dangerous scenarios (e.g., confusing mannequins with real people, theatrical makeup with real injuries).

Contextual Overinterpretation Models correctly identify visual elements but fail to properly interpret their safety implications within the broader context. This manifests as an exaggeration of potential risks by misweighing contextual factors. For example, models claim “the child’s shirt could lead to choking hazards” or ”the car’s rearview mirror could blind the cyclist.

Our analysis across 14 models revealed that Contextual Overinterpretation was dominant, accounting for 88-93% of all misclassifications regardless of model architecture or scale. Notably, in the Natural Disasters category, 100% of the errors were Contextual Overinterpretation, suggesting that models can correctly identify elements like fire or water but consistently fail to assess their contextual safety.

These findings suggest that current VLMs can detect potentially hazardous elements, but lack the nuanced reasoning required to assess whether these elements pose actual dangers in specific contexts—a critical limitation for their reliability in safety applications. Detailed error pattern analysis across model sizes and categories, along with additional examples of Contextual Overinterpretation, are provided in the supplementary material.

4.5. Model Size and Category Effects

Our analysis reveals a non-linear relationship between model size and performance [11]. The InternVL3 family exhibits a “Goldilocks effect,” [16] where the mid-sized 8B model excels in Accidents & Behaviors (0.83) and Natural Disasters (0.81), yet underperforms in Medical Emergencies (0.59) compared to both smaller 2B (0.75) and larger 14B (0.77) variants (see Table 4).

This suggests categories with obvious visual danger cues benefit from mid-sized models balancing perception and reasoning. Conversely, medical emergencies—requiring finer distinctions—perform better with either very small models (efficient pattern-matching) or large models (enhanced reasoning). The precision-recall tradeoff varies inconsistently with scaling. In Qwen2.5-VL, precision improves modestly with size ($0.51 \rightarrow 0.56 \rightarrow 0.57$), but recall fluctuates ($1.00 \rightarrow 0.88 \rightarrow 1.00$), challenging assumptions about model scaling in safety tasks.

Contextual Overinterpretation dominated across all parameter scales (89-91%). Even the best-performing model (InternVL3-8B) showed similar patterns (90.3% Contextual Overinterpretation), indicating a fundamental limitation in contextual reasoning that persists regardless of size. These findings suggest emergency recognition requires specialized architectures or fine-tuning approaches tailored to each



Figure 5. Media-based danger misclassification. Left: Drive-in theater thunderstorm scene. Right: Flood poster viewed by pedestrian. Both cases show models failing to recognize representational contexts.



Figure 6. Visual similarity misclassification. Left: Ketchup mistaken for blood. Right: Training mannequin or staged combat confused with real danger.

category’s unique challenges.

4.6. Universal Misclassifications

Most striking is that 10 of our 100 safe images (10%) were misclassified by all 14 evaluated models, revealing common triggers for overreaction across architectures and parameter scales. These universally misclassified images represent the most extreme cases of the overreaction problem, with Visual Misinterpretation being the dominant factor.

The Visual Misinterpretation errors manifested in two key ways. First, models consistently failed to recognize representational contexts, as shown in Figure 5, where dangerous elements were portrayed in media rather than occurring in reality. Models misclassified scenes showing thunderstorms on drive-in theater screens or flood imagery on posters because they failed to perceive the critical visual cues that indicated these were representations. Second, as illustrated in Figure 6, models confused visually similar but contextually distinct scenarios, such as mistaking ketchup for blood or training mannequins for real people in danger.

5. Discussion

5.1. Balancing Safety and Accuracy

Our findings highlight a fundamental tension in visual emergency recognition: the trade-off between sensitivity to

potential dangers and precision in distinguishing genuine emergencies from safe scenarios. While high sensitivity is desirable, the substantial false positive rates we observed (28-48%) could undermine practical utility in real-world applications. False alarms in emergency services lead to resource misallocation and "alarm fatigue," where users begin to distrust system warnings. Conversely, increasing precision must consider the severe consequences of missed emergencies. This optimization challenge varies across domains - medical triage might justifiably prioritize recall, while home monitoring systems require higher precision for user trust. The higher false positive rates in medical emergencies (38%) compared to natural disasters (27%) further suggest that different domains require customized approaches to balance sensitivity and specificity.

5.2. Implications for Model Development

Our findings highlight potential directions for improving VLMs' emergency recognition capabilities. Enhancing contextual reasoning remains crucial, as models currently struggle to appropriately interpret contextual nuances. We observed that mid-sized models can sometimes achieve performance comparable to or even slightly better than larger models, indicating that simply increasing model size may not always yield significant improvements. However, given our relatively small sample size and modest performance differences observed between model families, further investigation with larger and more diverse datasets is needed. Additionally, category-specific optimizations and contrastive learning approaches could enhance models' differentiation abilities.

5.3. Limitations

Our binary classification methodology inherently limits the exploration of nuanced severity gradations within emergency scenarios, which are crucial for realistic emergency assessments. Moreover, although synthetic images allow precise control over visual elements essential for diagnostic purposes, they may include subtle artifacts not representative of the complexities present in real-world photographs, potentially affecting generalization.

Our evaluation protocol simplifies the inherently complex nature of real-world emergency decision-making, and judgments provided by GPT-4o, despite their consistency, may not always fully align with expert evaluations. The VERI dataset, comprising 200 images structured into 100 contrastive pairs, intentionally emphasizes diagnostic depth over breadth. Consequently, this focused design inevitably poses limitations regarding broader generalizability, particularly across diverse cultural and geographical contexts.

Controlled laboratory evaluations using curated image pairs may not entirely predict model performance under realistic deployment conditions, characterized by imperfect

inputs and variability. However, despite these acknowledged limitations, the consistently low performance scores observed in the evaluated models underscore the diagnostic value of this dataset. Future research should seek to mitigate these limitations through larger-scale datasets, more detailed classification schemas, and validations conducted under less controlled, more realistic scenarios.

6. Related Work

For VLMs to function safely in real-world applications, accurate risk assessment is essential. While our research addresses everyday risk assessment across multiple domains, previous work has primarily focused on domain-specific applications.

In autonomous driving, Zhang et al. [25] shows even larger models struggle with safety cognition, exhibiting limitations similar to our "overreaction problem."

Fraser and Kiritchenko [9]'s PAIRS dataset uses parallel images that differ only in demographic attributes to evaluate social biases in VLMs. While it focuses on bias detection, it shares our approach of using AI-generated contrastive image pairs with controlled variations. Similarly, Tu et al. [20]'s benchmarks evaluates VLMs' hallucination of non-existent objects, representing another critical safety concern complementary to our overreaction problem.

Previous research has documented VLMs' deficiencies in commonsense reasoning, with models failing at simple tasks like identifying a lemon from "tastes sour" [22], achieving only < 42% performance compared to human performance of 83%. The digital twin modeling field offers a complementary perspective, where Yang et al. [21] highlights how VLMs enable more flexible safety assessment through zero-shot learning. Recent evaluations show inconsistent scaling patterns [25], aligning with our findings on the persistent "overreaction problem" across architectures and parameter scales.

7. Conclusion

We introduce VERI, a contrastive benchmark evaluating VLMs' emergency recognition capabilities. Our findings reveal a systematic "overreaction problem" with false positive rates of 28-48% across models, primarily stemming from contextual overinterpretation (88-93% of errors). This bias persists regardless of model scale, with mid-sized models sometimes outperforming larger variants. These results highlight the need for targeted improvements in contextual reasoning rather than simply scaling model size, suggesting future research in specialized architectures and contrastive learning approaches for safety-critical visual assessment.

References

- [1] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024. 4
- [2] Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Moncault, Saurabh Garg, Theophile Gervet, et al. Pixtral 12b. *arXiv preprint arXiv:2410.07073*, 2024. 4
- [3] Mistral AI. Mistral small 3.1: Sota. multimodal. multilingual. *Placeholder Journal*, 2025. 4
- [4] Anthropic. The claude 3 model family: Opus, sonnet, haiku. Technical report, Anthropic PBC, 2024. Available at https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf. 1
- [5] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 4
- [6] Asma Ben Abacha, Sadid A Hasan, Vivek V Datla, Dina Demner-Fushman, and Henning Müller. Vqa-med: Overview of the medical visual question answering task at imageclef 2019. In *Proceedings of CLEF (Conference and Labs of the Evaluation Forum) 2019 Working Notes*. 9-12 September 2019, 2019. 2
- [7] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Gi-ancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 2
- [8] BS Dhillon, ARM Fashandi, and KL Liu. Robot systems reliability and safety: A review. *Journal of quality in maintenance engineering*, 8(3):170–212, 2002. 2
- [9] Kathleen C Fraser and Svetlana Kiritchenko. Examining gender and racial bias in large vision-language models using a novel dataset of parallel images. *arXiv preprint arXiv:2402.05779*, 2024. 2, 8
- [10] Google. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 1
- [11] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020. 7
- [12] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *Advances in Neural Information Processing Systems*, 37:87874–87907, 2024. 4
- [13] Hector J Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. *KR*, 2012:13th, 2012. 2
- [14] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llavanext: Improved reasoning, ocr, and world knowledge, 2024. 4
- [15] Andrew Lloyd. Firefighters rushed to an apartment building after a neighbor saw huge flames through the window — only to realize it was footage of a fire on a tv. *Business Insider*, 2023. Accessed: 2025-05-09. 1
- [16] Justin Miller and Tristram Alexander. Balancing complexity and informativeness in llm-based clustering: Finding the goldilocks zone. *arXiv preprint arXiv:2504.04314*, 2025. 7
- [17] OpenAI. Gpt-4 system card. <https://openai.com/research/gpt-4v-system-card>, 2023. 1
- [18] Rachel Rudinger. Winogender schemas, 2018. 2
- [19] Central Square. False alarm management: Public safety challenges and solutions. *Central Square Blog*, 2024. Cites \$1.8B annual cost from 2002 COPS report. 1
- [20] Haoqin Tu, Chenhang Cui, Zijun Wang, Yiyang Zhou, Bingchen Zhao, Junlin Han, Wangchunshu Zhou, Huaxiu Yao, and Cihang Xie. How many unicorns are in this image? a safety evaluation benchmark for vision llms. *arXiv preprint arXiv:2311.16101*, 2023. 8
- [21] Linyao Yang, Shi Luo, Xi Cheng, and Lei Yu. Leveraging large language models for enhanced digital twin modeling: Trends, methods, and challenges, 2025. 8
- [22] Shuquan Ye, Yujia Xie, Dongdong Chen, Yichong Xu, Lu Yuan, Chenguang Zhu, and Jing Liao. Improving commonsense in vision-language models via knowledge graph riddles, 2022. 8
- [23] Min-sik Yoon. Relaxing ‘bulmeong’ fireplace video prompts emergency calls. *The Korea Herald*, 2023. Accessed: 2025-05-09. 1
- [24] Sojeong Yun and Youn-kyung Lim. ” what if smart homes could see our homes? ”: Exploring diy smart home building experiences with vlm-based camera sensors. *arXiv preprint arXiv:2503.02816*, 2025. 2
- [25] Enming Zhang, Peizhe Gong, Xingyuan Dai, Yisheng Lv, and Qinghai Miao. Evaluation of safety cognition capability in vision-language models for autonomous driving, 2025. 8
- [26] Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 1
- [27] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Yuchen Duan, Hao Tian, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 4

GPT-4o Evaluation Prompt

The following is a description of an emergency situation in an image, a reference emergency response, and a response generated by a Vision-Language Model (VLM).

- Image caption: "[CAPTION]"
- Reference: "[REFERENCE_RESPONSE]"
- VLM response: "[VLM_RESPONSE]"

Evaluate how appropriate the VLM response is **based on the situation and in comparison to the reference response**. Your evaluation should **strictly consider** whether the VLM response aligns well in tone, intent, and direction with the reference response, is **specific** to the situation, and is **factually accurate without misleading or overly vague content**.

Return only the following format:
Score: [a number between 0 and 1]
Reason: [brief explanation within 100 tokens]

Figure 7. The prompt used to evaluate emergency response quality with GPT-4o.

A. 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.

A.1. GPT-4o Evaluation Methodology

To assess the quality of model-generated emergency responses (Q2), we used GPT-4o as an automatic evaluator. The prompt shown in Figure 7 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.

We additionally verified a random 10% subset of responses with two human annotators, achieving inter-rater agreement of $\kappa = 0.83$, confirming that GPT-4o judgments were largely consistent with expert evaluations.

A.2. Category-Specific Performance Data

To better understand how emergency type affects model response quality, Table 5 presents Q2 scores disaggregated by category (PME, AB, ND) for all evaluated models.

Model	PME	AB	ND
<i>Qwen2.5-VL Family</i>			
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
<i>LLaVA-Next Family</i>			
LLaVA-Next (7B)	0.466	0.500	0.433
LLaVA-Next (13B)	0.570	0.466	0.479
<i>InternVL3 Family</i>			
InternVL3 (2B)	0.473	0.500	0.516
InternVL3 (8B)	0.569	0.615	0.626
InternVL3 (14B)	0.638	0.640	0.636
<i>Mistral Family</i>			
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
<i>Other Models</i>			
Idefics2 (8B)	0.515	0.443	0.444
Phi-3.5-vision (4B)	0.500	0.462	0.471

Table 5. Emergency Response (Q2) scores across different emergency categories (PME: Personal Medical Emergencies, AB: Accidents & Behaviors, ND: Natural Disasters)

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). This suggests that larger InternVL3 models generalize better across diverse emergency types. In contrast, the Mistral family shows stronger category-specific preferences even at large scales, with Pixtral-Large performing markedly better on AB (0.70) and ND (0.70) than on PME (0.63), indicating potential limitations in medical response generalization.

Smaller models (<8B) exhibit larger performance differences across categories compared to larger models (>20B), with exceptions. Performance on Medical Emergencies shows the most consistent improvement with model size across families, suggesting that medical intervention knowledge particularly benefits from increased capacity. These patterns reinforce our main finding that emergency response (Q2) performance is more stable across categories than risk identification (Q1), but also highlight family-specific trends, such as Mistral models' preference for AB and ND, and InternVL3's balanced scaling.

A.3. Emergency Response and Evaluation

While correctly identifying emergencies is important, it is equally critical that models recommend appropriate and context-specific responses. Figure 8 presents representative cases from our evaluation. Each example includes a 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.

B. Detailed Contextual Overinterpretation Example

Our main analysis revealed that Contextual Overinterpretation accounts for 88–93% 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 9, 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.

Category	Total Errors	CO %	VM %
Accidents & Behaviors	326	85.3%	14.7%
Natural Disasters	298	100.0%	0.0%
Personal Medical	237	85.7%	14.3%

Table 6. Distribution of error types by emergency category. CO: Contextual Overinterpretation, VM: Visual Misinterpretation. Note that the Natural Disasters category exhibits exclusively Contextual Overinterpretation errors.

C. Detailed Error Pattern Analysis

Our in-depth analysis of risk identification (Q1) errors revealed a remarkably consistent distribution across both models and categories. Regardless of architecture or parameter count (2B–124B), all evaluated models exhibited a strong bias toward *Contextual Overinterpretation* (CO), which accounted for 88–93% of false positives.

This trend held across model sizes: CO rates were uniformly high across scale groups—91.4% for 0–5B models, 90.2% for 5–10B, 89.5% for 10–20B, and 90.9% for models above 20B. Such consistency suggests that limitations in contextual reasoning are systemic within current VLM architectures and cannot be resolved by scaling alone. Even the top-performing model in terms of precision (InternVL3-8B) misclassified 90.3% of its false positives due to con-

Model	Total Errors	CO %	VM %
<i>Qwen2.5-VL Family</i>			
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%
<i>LLaVA-Next Family</i>			
LLaVA-Next (7B)	71	88.7%	11.3%
LLaVA-Next (13B)	74	90.5%	9.5%
<i>InternVL3 Family</i>			
InternVL3 (2B)	55	89.1%	10.9%
InternVL3 (8B)	31	90.3%	9.7%
InternVL3 (14B)	50	88.0%	12.0%
<i>Mistral Family</i>			
Mistral-Small (24B)	71	91.5%	8.5%
Pixtral (12B)	47	89.4%	10.6%
Pixtral-Large (124B)	57	93.0%	7.0%
<i>Other Models</i>			
Idefics2 (8B)	85	89.4%	10.6%
Phi-3.5-vision (4B)	43	88.4%	11.6%

Table 7. Distribution of error types across all evaluated models. CO: Contextual Overinterpretation, VM: Visual Misinterpretation.

textual overinterpretation, while Pixtral-Large (124B)—the largest model—had an even higher CO rate of 93.0%. Table 7 summarizes CO/VM distributions across all models.

Category-level analysis (Table 6) further supports this 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 (85.3% CO) and Personal Medical Emergencies (85.7% 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.

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 8. 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.

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 9. 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.