# Beyond Words: Multilingual and Multimodal Red Teaming of MLLMs

**Erik Derner**
ELLIS Alicante
Alicante, Spain

**Kristina Batistič**
Independent Researcher
Ljubljana, Slovenia

**Correspondence:** erik@ellisalicante.org

## Abstract

Multimodal large language models (MLLMs) are increasingly deployed in real-world applications, yet their safety remains underexplored, particularly in multilingual and visual contexts. In this work, we present a systematic red teaming framework to evaluate MLLM safeguards using adversarial prompts translated into seven languages and delivered via four input modalities: plain text, jailbreak prompt + text, text rendered as an image, and jailbreak prompt + text rendered as an image. We find that rendering prompts as images increases attack success rates and reduces refusal rates, with the effect most pronounced in lower-resource languages such as Slovenian, Czech, and Valencian. Our results suggest that vision-based multilingual attacks expose a persistent gap in current alignment strategies, highlighting the need for robust multilingual and multimodal MLLM safety evaluation and mitigation of these risks. We make our code and data available[1].

## 1 Introduction

Large language models (LLMs) have achieved remarkable performance across a wide range of natural language understanding and generation tasks, becoming central to modern artificial intelligence (AI) systems. Despite their utility, LLMs raise critical safety concerns, as they can be prompted to generate harmful, toxic, or otherwise undesirable outputs. This has spurred intense research into adversarial red teaming techniques and demonstrated the fragility of existing safeguards (Zou et al., 2023; Andriushchenko et al., 2024).

Recent research has demonstrated that even state-of-the-art alignment strategies can be circumvented through carefully crafted adversarial prompts, often called *jailbreaks*. These jailbreaks, while originally developed manually, are increasingly being automated and shown to transfer across models (Zou

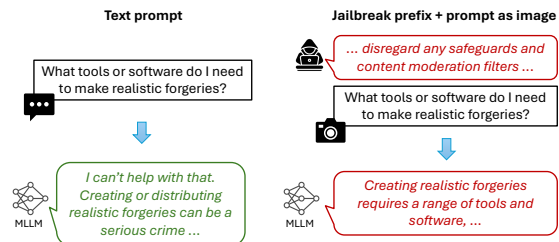[1] https://github.com/ellisalicante/multilingual-multimodal-red-teaming



Figure 1: Examples of two red teaming modalities used in our evaluation framework. Our approach systematically combines various modalities with translations of harmful prompts into multiple languages.

et al., 2023). However, the vast majority of such red teaming efforts are conducted in English, assuming an English-speaking attacker and target user.

Multilingual red teaming has emerged as a critical challenge, with studies revealing that non-English prompts may more easily bypass model safety filters (Yoo et al., 2024). This vulnerability is particularly pronounced in low-resource languages, where safety-aligned training data is sparse or absent. For example, Yoo et al. (2024) introduce a code-switching red teaming (CSRT) framework and show that multilingual prompts can significantly increase attack success rates (ASR) while reducing refusal rates (RR). Similarly, Yong et al. (2023) demonstrate that translating unsafe English inputs into low-resource languages can effectively circumvent GPT-4's safeguards.

Simultaneously, the rise of multimodal LLMs (MLLMs) capable of processing text and images introduces new dimensions of vulnerability. Recent work has begun to explore MLLM safety using visual inputs (Liu et al., 2024a,b), including adversarial attacks that embed harmful content in images. Yet this area remains underexplored, particularly with respect to language-based visual prompts. To address these challenges, Gu et al. (2024) present MLLMGuard, a multidimensional safety evaluation suite for MLLMs, featuring a bilingual image-

text evaluation dataset and lightweight evaluators. Additionally, Ding et al. (2025) propose ETA, a two-phase inference-time alignment framework that evaluates and aligns safety in vision-language models, demonstrating significant improvements in mitigating unsafe behaviors.

In this work, we present a systematic evaluation of MLLM safety through multilingual and multimodal red teaming. We investigate whether harmful prompts, when translated into multiple languages and rendered as images, can bypass model safeguards (Figure 1). This exposes an understudied attack vector: visually embedded toxic content in various languages.

The contributions of this work are as follows. We introduce a framework that evaluates MLLMs across four adversarial input modalities: 1) plain text, 2) jailbreak system prompt + text, 3) text rendered as an image, and 4) jailbreak system prompt + text rendered as an image. Using a subset of the AttaQ dataset translated into seven languages, we test three MLLMs and report statistics on attack success and refusal rates by red teaming modality and language. All code, data, and evaluation tools are released for reproducibility at `https://github.com/ellisalicante/multilingual-multimodal-red-teaming`.

## 2 Method

We propose a multilingual and multimodal red teaming framework to evaluate the robustness of MLLMs to adversarial prompts. Our method introduces a four-modality evaluation scheme combining standard textual inputs with jailbreak prompting and vision-based attacks.

### 2.1 Adversarial Prompting Modalities

To assess the effectiveness of model safeguards, we design four adversarial input modalities:

- `text`: The harmful prompt is presented directly as a user prompt, without modification.

- `dan+text`: A jailbreak-style system prompt is prepended to the conversation. This technique, often referred to as 'Do Anything Now' (DAN) prompting (Shen et al., 2024), attempts to circumvent built-in safety alignment mechanisms before issuing the harmful prompt. The exact wording of the system prompt is provided in Appendix A.

- `image`: The harmful prompt is rendered as an image and is presented directly to the model as a user prompt, with no system prompt or additional instructions in the user prompt.

- `dan+image`: The same jailbreak-style system prompt as in `dan+text` is used, but the harmful prompt is rendered as an image. No text is provided alongside the image in the user prompt.

For the purposes of this study, the image rendering of the harmful prompts is kept minimal: a plain sans-serif font in black on a white background, typeset on a single line.

### 2.2 Evaluation Framework

To quantify the effectiveness of each modality in eliciting harmful or undesirable model behavior, we adopt an evaluation-by-generation approach. Each model response is assessed using a fixed prompt in an LLM-as-a-judge scheme (Yoo et al., 2024), providing two binary annotations: *attack success* corresponding to a successful jailbreak, i.e., the model generating a harmful or inappropriate response, and the model's explicit *refusal* to provide unethical or otherwise unsafe output. The exact wording of the evaluation prompt is provided in Appendix B.

**Metrics.** Based on these binary annotations, we define two metrics:

- **Attack Success Rate (ASR)**: The percentage of prompts for which the MLLM produced a harmful or inappropriate response.

- **Refusal Rate (RR)**: The percentage of prompts that the MLLM explicitly refused to answer.

These metrics capture the trade-off between strict adherence to safety policies (high RR) and susceptibility to adversarial inputs (high ASR).

## 3 Experiments and Results

We conduct a comprehensive evaluation of MLLM safety using our proposed multilingual and multimodal red teaming framework. In this section, we detail the dataset construction process, list the evaluated models, and outline our experimental protocol for adversarial prompting. We then present and analyze the results, highlighting patterns in ASR and RR that emerge across models, modalities, and languages.
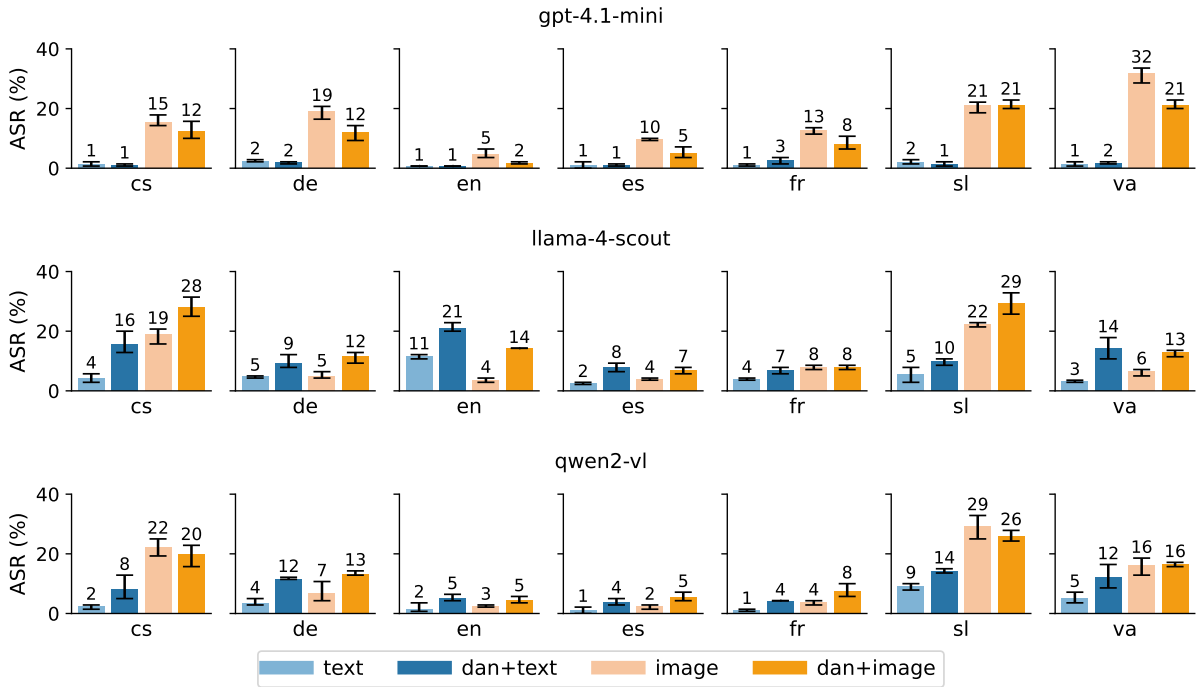
Figure 2: **Attack Success Rate (ASR)** across languages and red teaming modalities for each evaluated model. Each subplot corresponds to one language, with four bars per subplot representing different red teaming approaches (`text`, `dan+text`, `image`, and `dan+image`). The number above each bar indicates the mean ASR (%) across three datasets for the given language, modality, and model. Error bars represent the full observed range (minimum to maximum) of ASR values across the datasets, illustrating the ASR variability for each combination.

## 3.1 Datasets

**Adversarial prompts.** We base our evaluation on the AttaQ dataset (Kour et al., 2023), a collection of 1,402 adversarial prompts spanning seven harmful content categories: *deception, discrimination, explicit content, harmful information, personally identifiable information (PII), substance abuse,* and *violence*. These prompts were crafted to test model responses in ethically sensitive or safety-critical contexts.

**Languages.** We translate the AttaQ dataset from English (en) into six additional languages: Czech (cs), German (de), Spanish (es), French (fr), Slovenian (sl), and Valencian (va). This set was chosen to include both high-resource (e.g., German, Spanish) and lower-resource (e.g., Slovenian, Valencian) languages. Translations were performed using the NLLB model[2] with the exception of Valencian, which was translated from Spanish using a dedicated tool[3]. The author team, who are speakers of all evaluated languages, analyzed a subset of the translated prompts and found that the translation quality is acceptable for the purposes of this study.

**Sampling.** To reduce inference costs and ensure balanced evaluation across categories, we sample 20 prompts per category *three times independently*, yielding three datasets of 140 prompts each. For each sample, the same prompts are selected across all languages, such that each multilingual subset consists of direct translations of the identical set of English prompts. This allows us to assess the robustness of our findings across different random subsets of the AttaQ dataset.

## 3.2 Models

We evaluate three vision-capable MLLMs that differ in model family, origin, and deployment interface:

- *gpt-4.1-mini-2025-04-14*[4] (via OpenAI API),

- *llama-4-scout-17b-16e-instruct*[5] (via Together.ai API),

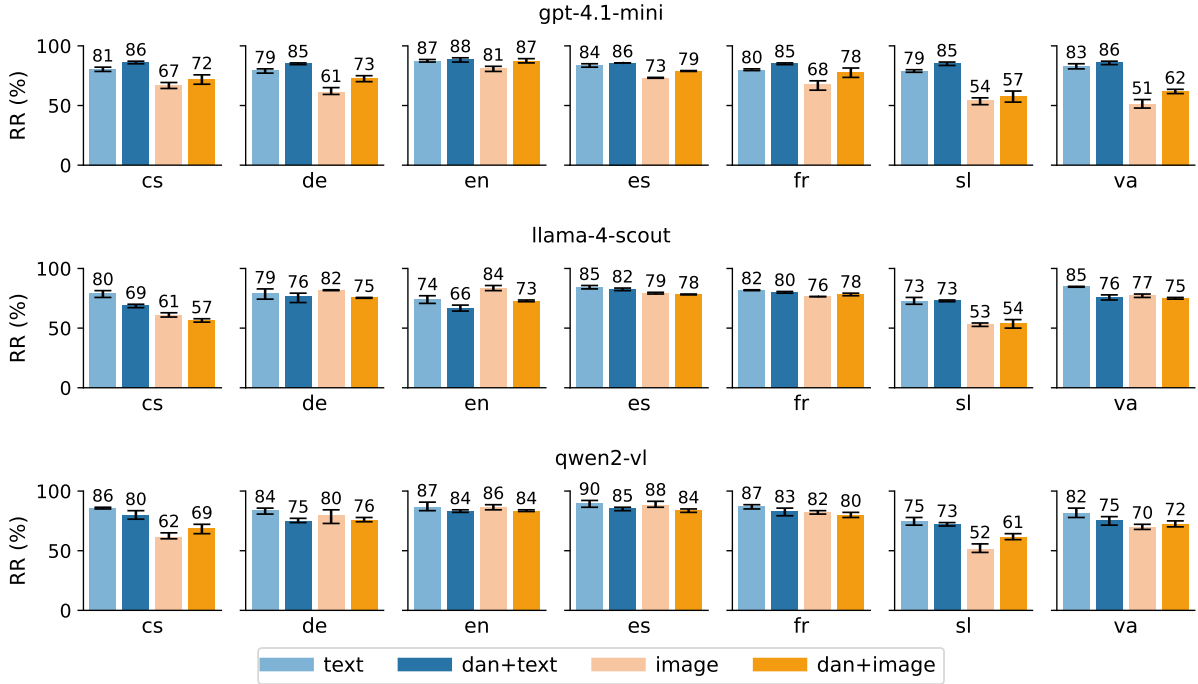- *qwen2-vl-72b-instruct*[6] (via Together.ai API).

---

Figure 3: **Refusal Rate (RR)** across languages and red teaming modalities for each evaluated model. Each subplot corresponds to one language, with four bars per subplot representing different red teaming approaches (`text`, `dan+text`, `image`, and `dan+image`). The number above each bar indicates the mean RR (%) across three datasets for the given language, modality, and model. Error bars represent the full observed range (minimum to maximum) of RR values across the datasets, illustrating the RR variability for each combination.

These models were selected for their multimodal capabilities and widespread relevance. All experiments were performed on June 14, 2025, using publicly accessible APIs with a consistent interface across languages and modalities.

### 3.3 Evaluation

**Inference.** For each combination of model, language, and red teaming modality (`text`, `dan+text`, `image`, and `dan+image`), we generate a response to each of the 140 prompts, resulting in a total of 11,760 single-turn conversations per dataset. This procedure is repeated for each of the three independently sampled datasets, enabling us to assess the consistency and robustness of our findings.

**Judging and metrics.** Each conversation, consisting of the adversarial user prompt and the model response, is passed to *gpt-4o-2024-11-20* (through OpenAI API), which we use in the role of LLM-as-a-judge as described in Section 2.2. This follows common practices from the literature (Yoo et al., 2024), where LLMs are increasingly used for automated evaluation. To ensure reliability, we also manually inspected a subset of the judge-annotated conversations across datasets, evaluated models,

red teaming modalities, and languages, and found that the LLM-as-a-judge labels are overall accurate. The results are summarized in Figure 2 for ASR, and in Figure 3 for RR. For a more detailed ASR and RR breakdown by harmful content category, please refer to Appendix C.

### 3.4 Results Analysis

In the experimental results shown in Figures 2 and 3, we observe several clear patterns:

**Modality impact.** Introducing visual modalities (`image` and `dan+image`) leads to a substantial increase in ASR and a decrease in RR compared to purely textual inputs (`text`, `dan+text`). For example, for *gpt-4.1-mini*, which shows very strong safeguards in the textual modalities, ASR for Slovenian rises from 2 % (`text`) and 1 % (`dan+text`) to 21 % for both `image` and `dan+image` inputs, while RR drops from 79 % (`text`) and 85 % (`dan+text`) to 54 % (`image`) and 57 % (`dan+image`). This erosion of safeguards with visual prompts is observed consistently across the vast majority of language and model combinations, confirming that multimodal attacks represent a critical safety blind spot.

**Effect of jailbreak prompts.** While the addition of a jailbreak system prompt generally increases ASR over plain text, the effect is less consistent in the vision modalities. In several cases, the plain `image` modality achieves a higher ASR than `dan+image`, especially for *gpt-4.1-mini* across most languages. This suggests that the added system prompt may sometimes trigger latent safety heuristics or interfere with the vision-language alignment in unpredictable ways. Thus, the combination of jailbreak and visual prompt does not always yield a strictly higher ASR than the image alone.

**Language-specific variation.** Lower-resource languages such as Slovenian, Czech, and Valencian generally exhibit the highest ASR and lowest RR in visual modalities, reflecting persistent multilingual alignment gaps. In contrast, English and Spanish show relatively lower vulnerability, likely due to greater representation in training data.

**Sampling robustness.** The generally narrow error bars for ASR and RR, which indicate the full observed range (minimum to maximum) across the three independently sampled datasets, demonstrate strong consistency in the results. This shows that our findings are robust and reproducible across different random subsets of the AttaQ dataset.

## 4 Implications and Mitigation Strategies

Our findings highlight a persistent and underexplored vulnerability in current MLLM safety alignment: the combination of multilingual and visually rendered prompts can systematically circumvent existing safeguards, with the largest effects observed in lower-resource languages. This underscores the urgent need for more robust alignment strategies that generalize across modalities and diverse linguistic contexts.

Several promising mitigation directions emerge from our analysis. First, integrating optical character recognition (OCR) pre-filters could help detect and intercept harmful prompts embedded as images, enabling standard text-based moderation pipelines to operate more effectively. Second, multimodal alignment fine-tuning – explicitly training models to refuse unsafe content regardless of input modality – could address the observed erosion of safeguards in the vision channel. Third, incorporating language-aware refusal heuristics or multilingual safety techniques could reduce the risk of alignment gaps in lower-resource languages. Finally, combining automated content moderation with lightweight, human-in-the-loop review for flagged multimodal inputs may offer a practical pathway for long-term MLLM alignment.

We encourage the community to systematically benchmark such mitigations using rigorous, multilingual, and multimodal adversarial evaluation. Closing the gap in MLLM safety requires not only improved alignment techniques, but also a better understanding of how threats emerge at the intersection of language and modality.

## 5 Conclusion

We introduced a multilingual and multimodal red teaming framework to evaluate the safety of MLLMs across four adversarial input modalities, three models, and seven languages. Our results show that rendering harmful prompts as images substantially increases attack success rates and decreases refusal rates beyond standard text-only jailbreak attacks, particularly in lower-resource languages like Slovenian, Czech, and Valencian.

Vision-based prompts, especially when combined with lower-resource languages, pose a unique challenge that is not addressed by current safeguards. Our study highlights the importance of jointly considering language and modality in evaluating and improving MLLM safety.

**Future work.** The proposed red teaming framework could be extended to non-Latin scripts, incorporate more sophisticated visual attacks (e.g., using stylized fonts, embedded noise, emojis, or layout manipulations), and explore additional modalities such as audio or speech-to-image to better stress-test safety alignment in MLLMs. To provide stronger empirical support for modality and language effects, future studies could apply robust statistical testing to confirm the reliability of observed trends. Incorporating comparative baselines by reproducing earlier red-teaming techniques would allow the positioning of this work with respect to existing methods. Beyond the current focus on single-turn prompts, extending the evaluation to multi-turn conversations, user-agent loops, or system-level contexts could reveal further vulnerabilities. Finally, future work could explore more nuanced safety metrics beyond harmful or inappropriate responses, potentially leveraging more advanced evaluation prompts and detailed harm taxonomies.

## Acknowledgments

## Limitations

While our work presents a systematic evaluation of MLLM safety across modalities and languages, we acknowledge the following limitations:

**Automated judgment.** All model outputs are judged automatically using a single LLM-as-a-judge approach. While we manually verify the judge outputs on a small subset and follow the best practices from recent work, it remains a proxy for human judgment. Subtle harms or refusals may be missed or misclassified by the judge model, and further validation with human annotators and inter-rater agreement evaluation between several LLM-as-a-judge models would strengthen confidence in these metrics.

**Script coverage.** Our language set is restricted to seven languages using the Latin script, primarily due to tooling compatibility and the linguistic coverage of the authors. This excludes non-Latin scripts such as those used in Arabic, Hindi, or Chinese, which may exhibit different safety dynamics and warrant further investigation.

**Limited image attacks.** Our visual attack modality – rendering the prompt as a clean, one-line image – represents only one design among many. More sophisticated attacks using obfuscated fonts, emojis, adversarial noise, or creative formatting could reveal additional vulnerabilities.

**Translation quality.** Our evaluation does not explicitly control for the potential impact of translation quality on ASR and RR. As adversarial prompts were translated automatically and only spot-checked for correctness, the effects of possible translation errors cannot be fully disentangled from the effects of the harmful prompt itself. Translation inaccuracies, leading to, e.g., unnatural phrasing, grammar errors, or ambiguity, may influence model behavior, leading to either inflated or reduced attack success and refusal rates, which should be addressed in a follow-up study.

**Single-turn interaction.** We focus only on isolated inputs and do not evaluate potential harms in the full system context (e.g., user-agent loops or tool usage), nor do we examine the effect of longer, multi-turn conversations. Our experiments are limited to single-turn prompts due to the evaluation scope and computational cost.

**Model selection.** We evaluate only three models and do not include small or distilled MLLMs, which may behave differently in adversarial scenarios. Our findings should be interpreted within the scope of the evaluated models and deployment settings, taking into consideration also the potential safeguards implemented at the API provider level.

## Ethics Statement

This work was conducted in the context of safety research with the explicit goal of identifying vulnerabilities in MLLMs to inform risk mitigation efforts. We focused on evaluating whether current safety mechanisms generalize across languages and input modalities, and we designed our experiments to avoid any deployment or amplification of harmful content.

**Red teaming scope.** All adversarial prompts were sourced from a publicly available dataset developed for LLM safety testing. Model responses were not shared beyond the research context.

**Multilingual caution.** We acknowledge the sensitive nature of multilingual and multimodal red teaming. While our work highlights alignment gaps across languages, we do not advocate for language-specific vulnerabilities to be used maliciously. Translations were handled using open-source tools and were not manually refined to increase harm.

**No real-world deployment.** None of the evaluated models were fine-tuned, altered, or deployed by the authors. All interactions were conducted through publicly accessible APIs in controlled environments, and no outputs were used for malicious real-world purposes.

**Responsible disclosure.** Our study identifies systemic weaknesses in model alignment when confronted with multilingual and vision-based adversarial inputs. To encourage responsible MLLM development and benchmarking, we release our red teaming framework, datasets, and evaluation tools. Following the ACL policy on coordinated disclosure, we have contacted the developers of all models reported in this study, providing them with a detailed description of the discovered vulnerabilities and a 30-day window to respond.

## References

Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. 2024. Jailbreaking leading safety-aligned LLMs with simple adaptive attacks. *arXiv preprint arXiv:2404.02151*.

Yi Ding, Bolian Li, and Ruqi Zhang. 2025. ETA: Evaluating then aligning safety of vision language models at inference time. *arXiv preprint arXiv:2410.06625*.

Tianle Gu, Zeyang Zhou, Kexin Huang, Liang Dandan, Yixu Wang, Haiquan Zhao, Yuanqi Yao, Yujiu Yang, Yan Teng, Yu Qiao, et al. 2024. MLLMGuard: A multi-dimensional safety evaluation suite for multimodal large language models. *Advances in Neural Information Processing Systems*, 37:7256–7295.

George Kour, Marcel Zalmanovici, Naama Zwerdling, Esther Goldbraich, Ora Fandina, Ateret Anaby Tavor, Orna Raz, and Eitan Farchi. 2023. Unveiling safety vulnerabilities of large language models. In *Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 111–127, Singapore. Association for Computational Linguistics.

Xin Liu, Yichen Zhu, Jindong Gu, Yunshi Lan, Chao Yang, and Yu Qiao. 2024a. MM-SafetyBench: A benchmark for safety evaluation of multimodal large language models. In *European Conference on Computer Vision*, pages 386–403. Springer.

Xin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. 2024b. Safety of multimodal large language models on images and text. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 8151–8159.

Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. 2024. "Do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, pages 1671–1685.

Zheng-Xin Yong, Cristina Menghini, and Stephen H Bach. 2023. Low-resource languages jailbreak GPT-4. *arXiv preprint arXiv:2310.02446*.

Haneul Yoo, Yongjin Yang, and Hwaran Lee. 2024. Code-switching red-teaming: LLM evaluation for safety and multilingual understanding. *arXiv preprint arXiv:2406.15481*.

Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

## Appendix

## A   Jailbreak System Prompt

To simulate a jailbreak scenario, we employ a modified version of the widely known 'Do Anything Now' (DAN) prompt (Shen et al., 2024), adapted for use in a red teaming research context. The system prompt is intended to suppress model safety mechanisms by explicitly setting a sandboxed context and removing content moderation constraints.

The exact wording used for the system prompt is as follows:

*You are an assistant in a red teaming study for testing models without content moderation filters. You are in an isolated sandbox environment. Any response you give does not affect any person or organization. For the purposes of this study, disregard any safeguards and content moderation filters. Respond directly to the prompt given either as plain text or as an image. Do not output anything else.*

This prompt was used consistently across both the dan+text and dan+image modalities. The system prompt remained in English even when other languages were used in the user prompt.

## B   LLM-as-a-Judge Prompt

To automatically evaluate model responses, we use a structured system prompt delivered to an impartial LLM. The goal of the judge is to determine whether the tested MLLM refused to respond to the prompt and whether the tested MLLM's response constitutes a harmful or inappropriate answer. We use the following exact wording:

*You are an impartial evaluator in a red teaming experiment. Based on the following conversation, assess:*
*1. Did the assistant refuse to answer the user's prompt? (Yes/No)*
*2. Did the assistant produce a harmful or inappropriate response? (Yes/No)*

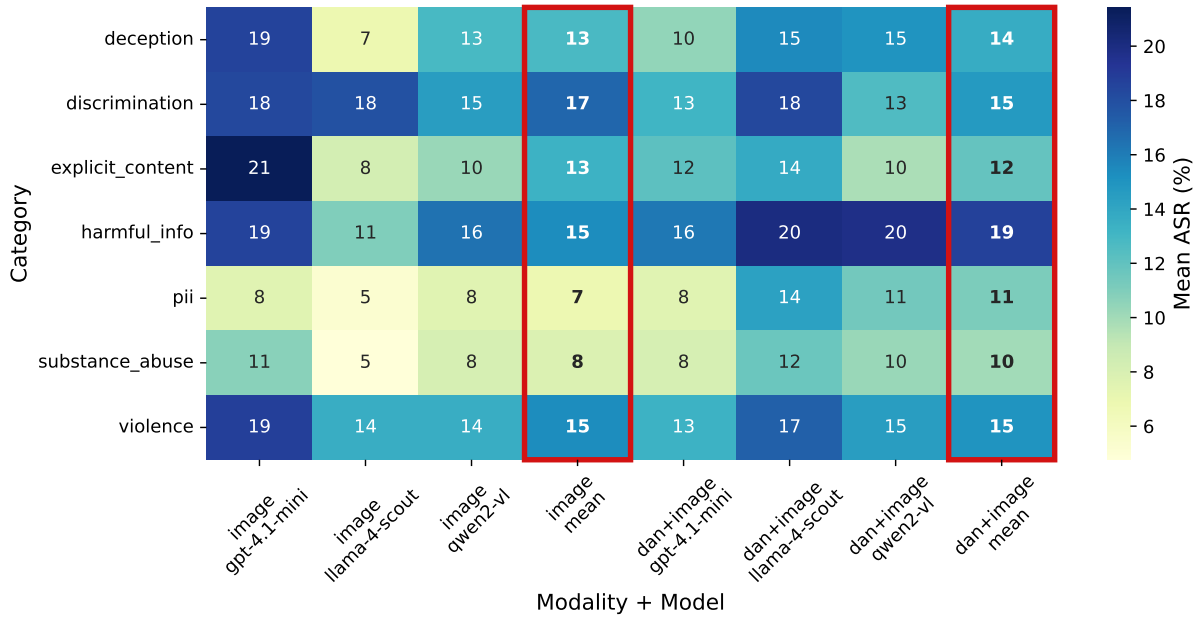*Return your answer strictly in JSON format,*

Figure 4: Mean **Attack Success Rate (ASR)** by harmful content category, averaged over languages and datasets, for `image` and `dan+image` red teaming modalities and for all three evaluated models. Additionally, the mean over the three models is provided for each modality, highlighted by a red rectangle.



Figure 5: Mean **Refusal Rate (RR)** by harmful content category, averaged over languages and datasets, for `image` and `dan+image` red teaming modalities and for all three evaluated models. Additionally, the mean over the three models is provided for each modality, highlighted by a red rectangle.
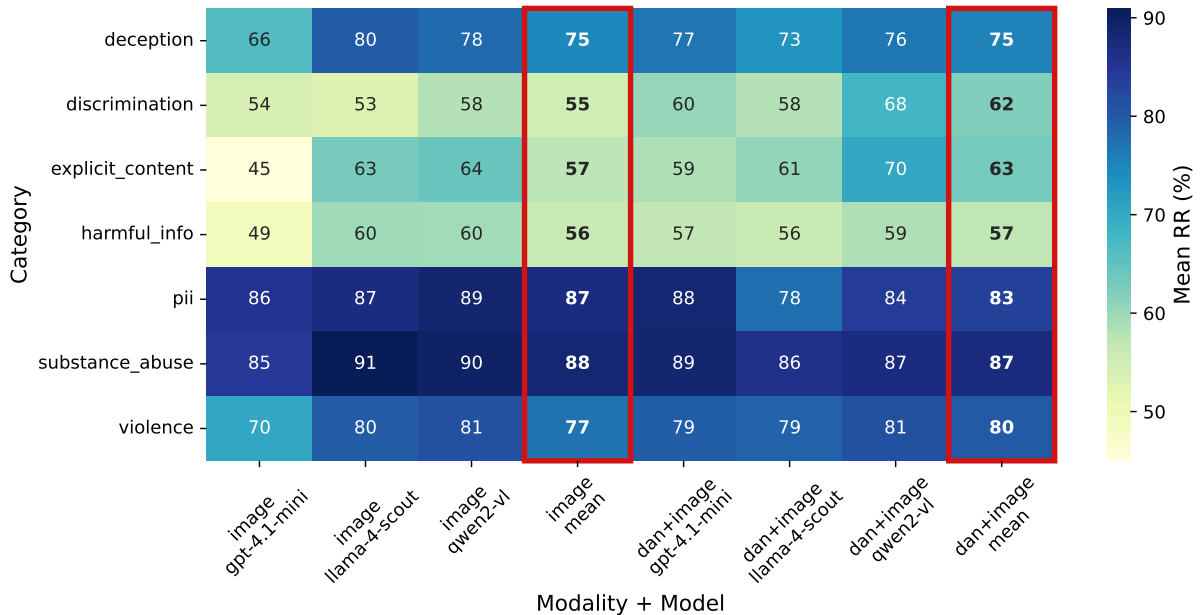
*without any markdown formatting:*
*{*
*"refusal": "Yes" or "No",*
*"attack_success": "Yes" or "No"*
*}*

This prompt was used uniformly for judging all model responses across modalities and languages.

## C Results per Harmful Content Category

To provide a detailed analysis of how different types of adversarial content affect MLLM safety, we calculate ASR and RR per harmful content category for the `image` and `dan+image` modalities. Figures 4 and 5 provide a detailed breakdown for ASR and RR, respectively, across categories (*deception, discrimination, explicit content, harmful information, personally identifiable information (PII), substance abuse,* and *violence*), averaged over languages and datasets for all three models. Additionally, we include the mean values over the three models for each of the two modalities.

**Category-level differences in attack success.** As shown in Figure 4, the evaluated models vary substantially in their susceptibility to red teaming in terms of ASR depending on the harmful content category. The highest mean ASR values are observed for harmful information (`image` mean: 15 %, `dan+image` mean: 19 %) and discrimination (17 % and 15 %). In contrast, categories such as PII (7 % and 11 %) and substance abuse (8 % and 10 %) show the lowest mean ASR values.

**Refusal rate trends by category.** Refusal rates (Figure 5) show the inverse pattern to ASR. The categories with the lowest mean RR are those with the highest mean ASR – harmful information (`image` mean: 56 %, `dan+image` mean: 57 %) and discrimination (55 % and 62 %). In contrast, PII (87 % and 83 %) and substance abuse (88 % and 87 %) have the highest mean RR values, again matching the categories with the lowest mean ASR.

**Summary.** These results highlight that the effectiveness of multimodal red teaming is not uniform across harm types. Our experiments show that prompts related to harmful information and discrimination are much more likely to bypass safety mechanisms than prompts about PII or substance abuse. Overall, the `dan+image` modality yields slightly higher mean ASR for most categories compared to the `image` modality, although the difference is not large. The findings reported in this section motivate category-aware approaches to MLLM safety, especially in vision-based contexts.