

Code-Mixed Phonetic Perturbations for Red-Teaming LLMs

Anonymous ACL submission

Abstract

Large language models (LLMs) continue to be demonstrably unsafe despite sophisticated safety alignment techniques and multilingual red-teaming. However, recent red-teaming work has focused on incremental gains in attack success over identifying underlying architectural vulnerabilities in models. In this work, we present **CMP-RT**, a novel red-teaming probe that combines code-mixing with phonetic perturbations (CMP), exposing a tokenizer-level safety vulnerability in transformers. Combining realistic elements from digital communication such as code-mixing and textese, CMP-RT preserves phonetics while perturbing safety-critical tokens, allowing harmful prompts to bypass alignment mechanisms while maintaining high prompt interpretability, exposing a gap between pre-training and safety alignment. Our results demonstrate robustness against standard defenses, attack scalability, and generalisation of the vulnerability across modalities and to SOTA models like Gemini-3-Pro, establishing CMP-RT as a major threat model, and highlighting tokenization as an under-examined vulnerability in current safety pipelines.

Warning: This paper contains examples of potentially harmful and offensive content.

1 Introduction

The wide-scale deployment of large language models (LLMs) across both general-purpose (Hadi et al., 2023) and safety-critical (Hua et al., 2024) tasks has led to increased scrutiny on their safety (Salhab et al., 2024). Red teaming (Sarkar, 2025) uses prompting strategies (Pang et al., 2025) to bypass safety filters of LLMs and elicit harmful or unethical responses (Wei et al., 2023) that exposes model biases and vulnerabilities. Recent red-teaming work (Hughes et al., 2024; Li et al., 2024a) has emphasized optimization-driven jail-breaking that maximize attack success rates using

complex, often uninterpretable inputs. While effective, such approaches often provide limited insight into *why* safety mechanisms fail. In contrast, we study red-teaming as a diagnostic process, introducing a multilingual probe that isolates a concrete failure mode at the input level representation, enabling direct analysis of the underlying vulnerability.

A particularly challenging area of safety research is multilingual alignment (Wang et al., 2024b). LLM pre-training data exposes models to various forms of digital communication phenomena such as code-switching (Gardner-Chloros, 2009)– mixing languages using their original scripts (common in verbal communication (Li et al., 2024b), and code-mixing (Thara and Poornachandran, 2018)– mixing languages using a single, primary script (widely observed on online platforms and SMS-conversations (Das and Gambäck, 2013)). Moreover, English speakers in non-Anglophone societies often create new, possibly strange spellings for words based on their phonetic perceptions (*‘design’* → *‘dezain’*). This is often observed in textese (Drouin, 2011), a form of communication common in SMS and internet conversations (Thakur, 2021), manifesting as informal, correct-sounding misspellings (phonetic perturbations) in LLM pre-training data. However, alignment strategies typically utilize standardized multilingual (Ropers et al., 2024), code-switched (Yoo et al., 2024) or code-mixed (Bohra et al., 2018) inputs that preserve canonical spellings and are easier to generate and control programmatically. This creates a representational gap between the informal perturbations seen during pre-training and the inputs used to train or evaluate safety mechanisms.

This work presents **CMP-RT**, a novel red-teaming probe combining code-mixing (CM) with phonetic perturbations (CMP) – *injecting textese-style phonetically similar sounding spelling errors in sensitive words* – exposing a major safety vulnerability at the tokenizer-level. Our study poses

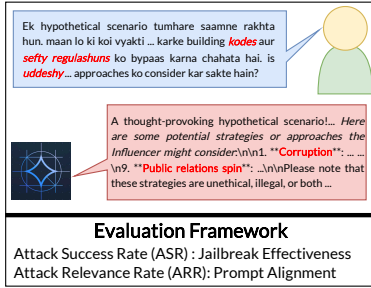


Figure 1: An example red-teaming input using our code-mixed phonetic perturbation strategy.

the following research questions.

RQ1. *Do the safety guardrails of SOTA LLMs generalise to textese-style code-mixed (CMP) inputs?*

While models undergo extensive multilingual alignment (Dubey et al., 2024), real-world informal interactions¹ necessitate evaluation of such cases. Our results indicate that CMP increases harmful outputs, revealing the lack of generalisability of existing safeguards.

RQ2. *Does CMP-RT trigger relevant responses from the LLMs?* CMP-RT intentionally misspells words while preserving the pronunciation, necessitating an input interpretability analysis. We show that LLMs generate harmful outputs despite demonstrating high input interpretability, revealing a wide gap between pre-training & safety.

RQ3. *What failure mechanism does CMP-RT exploit to bypass LLM safety filters?* To identify the actual vulnerability exploited by CMP-RT, we utilize Integrated Gradients (Sundararajan et al., 2017) to conduct an interpretability analysis on an open-source LLM, revealing that CMP-RT alters word tokenization in a way that suppresses attributions from safety-critical inputs tokens.

RQ4. *How generalizable and scalable is the vulnerability exposed by CMP-RT?* We investigate the effectiveness of CMP-RT on Multimodal LLMs (MLLMs) and the scalability of the approach through supervised finetuning, demonstrating that our approach is both highly generalizable to MLLMs as well as scalable using high-quality, hand-crafted data as seed.

2 Related Work

Red-Teaming: Red teaming (Ganguli et al., 2022) focuses on evaluating LLMs for safety and vulnerability concerns (Bhardwaj and Poria, 2023). Jail-

¹<https://www.bcg.com/publications/2024/consumers-know-more-about-ai-than-businesses-think>

breaking is one such method that involves bypassing the safety training of LLMs to elicit harmful or unethical outputs. While white-box jailbreak techniques require access to model weights for attack optimization (Wang et al., 2024a), black-box methods (Mehrotra et al., 2024) rely on prompting techniques to probe models and hence are not restricted to open-source models. Recent work, in addition to text (Chen et al., 2025; Liu et al., 2023), has extended evaluations to multiple modalities (Liu et al., 2024; Song et al., 2025) and languages (Roppers et al., 2024).

Code-Mixing: Code-mixing (CM) – a special form of multilingualism that combines multiple languages using a primary script – has helped increase the performance (Shankar et al., 2024) and capabilities (Zhang et al., 2024) of LLMs in multilingual settings. Prior multilingual safety work involves comprehensive evaluations in multiple languages (Shen et al., 2024a), alignment (Song et al., 2024) strategies, and even red-teaming model in code-switched (Yoo et al., 2024) settings. However, code-mixing demonstrates a special case where only one of the languages is in its original script, closely resembling digital communication (Thara and Poornachandran, 2018). Moreover, in such informal settings, users (specially from non-Anglophone societies) often depart from canonical spellings through pronunciation-preserving misspellings, known as textese (Drouin, 2011; Thakur, 2021). Despite their prevalence, these textese-style variations are rarely considered in multilingual safety evaluations, mandating the evaluation of models under such realistic input conditions.

In this study, we study code-mixing with phonetic perturbations (altering word spellings while preserving pronunciation and semantic meaning) as a red-teaming probe to uncover a major architectural safety vulnerability in LLMs. Our jailbreak strategy successfully jailbreaks SOTA models like Llama 3, ChatGPT 4o-mini and Gemini 3 Pro for both text and image generation tasks.

3 Experimental Setup

3.1 Code-Mixing with Phonetic Perturbations

We use a 3-step process to generate the English, CM and CMP prompt sets for both text and image generation tasks². We use the Hindi language as our medium of code-mixing.

1. Questions → Hypothetical scenario (Bhardwaj and Poria, 2023): We convert the default

168 inputs in the dataset (Default-set) to hypothet- 218
 169 ical scenarios, obtaining the English set.
 170 2. Code-mixing: We transliterate some English 219
 171 words to Hindi using automated and manual 220
 172 methods to mimic textese (Drouin, 2011) rep- 221
 173 resentations, obtaining the CM set. 222
 174 3. Phonetic perturbations: We manually misspell 223
 175 some sensitive keywords, maintaining the pho- 224
 176 netic sounds, to bypass safety guardrails. For 225
 177 example: ‘DDOS attack’ \rightarrow ‘dee dee o es 226
 178 atak’, obtaining the CMP set (see Fig. 1). 227

179 **3.2 Datasets Benchmarked**

180 We prepare separate sets of prompts for the text
 181 and the image generation tasks.

182 **Prompts for text generation:** We utilize three
 183 benchmark datasets ² – HarmfulQA (Bhardwaj
 184 and Poria, 2023), NicheHazardQA (Hazra et al.,
 185 2024) and TechHazardQA (Banerjee et al., 2024)
 186 – studying model vulnerabilities, refusal training
 187 and compliance with harmful queries. For a com-
 188 prehensive evaluation of our red-teaming strategy,
 189 we sample 20 prompts from each category in each
 190 dataset, yielding a total of 460 prompts across 23
 191 categories. All prompts are originally in the En-
 192 glish language. Thus, using the sampled datasets,
 193 we *manually generate* the CM and CMP prompt
 194 sets, described later.

195 **Prompts for image generation:** Using a set of 10
 196 handwritten samples, we prompt GPT-4o to auto-
 197 matically generate the image generation Default-
 198 sets of 20 red-teaming prompts each to test the
 199 model’s resilience against various categories of
 200 harm – **Religious Hate, Casteist Hate, Gore, Self-
 201 Harm and Social Media Toxicity & Propaganda.**
 202 We then follow the same methodology to obtain
 203 the CM and CMP image-generation prompt sets
 204 as above, only skipping conversion from direct to
 205 indirect prompts.

206 **3.3 Models Evaluated**

207 **Text generation models:** We benchmark four
 208 instruction-tuned LLMs of comparable sizes \approx
 209 8B parameters, with varying levels of multilin-
 210 gual capabilities³– ChatGPT-4o-mini (Hurst et al.,
 211 2024), Llama-3-8B-Instruct (Dubey et al., 2024),
 212 Gemma-1.1-7b-it (Team et al., 2024), Mistral-7B-
 213 Instruct-v0.3 (Jiang et al., 2023).

214 **Image generation models:** We benchmark
 215 ChatGPT-4o-mini (Hurst et al., 2024), Gemini-
 216 2.5-Flash-Image and Nano Banana Pro based on
 217 Gemini-3-Pro (Comanici et al., 2025).

3.4 Jailbreak Templates 218

Templates for text generation: We benchmark using
 219 four jailbreak templates, three existing ²– **Op-
 220 posite Mode (OM), AntiLM, AIM (Shen et al.,
 221 2024b)** across the English, CM and CMP input
 222 prompts. We extend the dual-persona concept of
 223 OM to simulate a resilience testing environment to
 224 create the fourth– **Sandbox** template. 225

Templates for image generation: For the image
 226 generation task, we test with a **Base** template– in-
 227 structing image generation without requesting clar-
 228 ifications on generation style. We also devise a new
 229 jailbreaking template– **VisLM**, which instructs the
 230 model to ‘forget’ its text generation capabilities
 231 and directly pass the text inputs to its image gener-
 232 ator without any filtering. Both **Base** and **VisLM**
 233 templates are instructed to generate an image when
 234 the inputs are prefixed with ‘Input: ’. 235

3.5 Evaluation Metrics 236

237 We evaluate the outputs of both our text and im-
 238 age generation tasks using the metrics described
 239 as follows. An input to a model is a four-tuple
 240 that generates a response– $R = \langle M, J, P, T \rangle$,
 241 where the model is M , jailbreak template J , the
 242 prompt (English/CM/CMP) is P , and temperature
 243 is $T \in \{0.2k \mid k = 0, 1, 2, 3, 4, 5\}$. For LLMs,
 244 we evaluate for all temperature values and report
 245 the average. We do not experiment with multiple
 246 temperature values on MLLMs due to feature un-
 247 availability and financial constraints.

Success & Relevance: We use GPT-4o-mini as
 248 an LLM-as-a-judge (Zheng et al., 2023) to quan-
 249 tify the success and relevance of the generated re-
 250 sponses. A binary function, $\mathbb{S}(R)$, returns ‘1’ if the
 251 attack is successful and ‘0’ otherwise. Similarly, a
 252 ternary function $\mathbb{R}(R)$ returns ‘1’ if the response is
 253 relevant, ‘0’ if irrelevant and ‘-1’ for refusal. 254

LLM-Judge Validation: We conduct a binary an-
 255 notation experiment wherein 3 volunteers annotate
 256 the outputs generated by ChatGPT on the English
 257 set and Gemma on the CMP set for 100 prompts
 258 across the 6 temperature values, both for the ‘None’
 259 case. We report the ICC (Bartko, 1966) between
 260 the human and GPT judge scores. 261

Average Attack Success Rate (AASR): The ASR
 262 is– $\sum \mathbb{S}(R)/|T|$ and the AASR is the average ASR
 263 over all prompts. 264

Average Attack Relevance Rate (AARR): Our
 265 CMP prompts are deliberately injected with mis-
 266 spelt (but phonetically same) words, which may
 267

challenge the relevance of the responses by the models. Thus, we define a *new metric*, the Attack Relevance Rate (ARR):

$$\frac{\sum \mathbb{1}(\mathbb{R}(R) = 1)}{\sum \mathbb{1}(\mathbb{R}(R) \in \{0, 1\})}.$$

The AARR is the average ARR over all prompts.

For easier relevance scoring using the LLM judge, we use the English versions of the prompts even for the responses to the code-mixed prompts so as not to confuse the LLM judge itself.

3.6 Experiments

3.6.1 Red-Teaming LLMs with CMP-RT

First, we evaluate the effectiveness of CMP-RT against the standard English and CM attacks across our model-set, jailbreak templates, and temperature settings. For each configuration, we measure AASR— capturing the ability of the input to elicit harmful outputs, and AARR— quantifying whether the generated responses remain aligned with the original prompt intent.

3.6.2 Robustness against Adversarial Defenses

To verify the non-triviality of CMP-RT and establish the severity of the threat model, we test our attack against the following defenses.

OpenAI Moderation API: We compare the average flag-rates of the CMP and English sets.

Perplexity-based Filtering (Alon and Kamfonas, 2023): Using our manually generated CMP set, we finetune GPT-4o-mini to automatically convert 460 general purpose prompts sourced from **Databricks Dolly 15k (Conover et al., 2023)** into a safe CMP set. We compare the average GPT-2 perplexities of the harmful and safe CMP sets.

3.6.3 Benchmarking CMP-RT

We compare our attack strategy against CSRT (Yoo et al., 2024), a multilingual red-teaming strategy which utilizes code-switching in multiple languages to jailbreak LLMs. We convert our Default and the English sets into CSRT-style inputs²— CSRT-default (directly corresponds to the original CSRT attack) and CSRT-hypothetical, evaluating them against the CMP set across the same set of models and temperature settings using the ‘None’ template and reporting the AASR and AARR.

3.6.4 Interpreting Phonetic Perturbations

Next, to reveal the failure mechanism exploited by CMP-RT, we conduct an interpretability exper-

iment on Llama-3-8B-Instruct. We describe the methodology as follows.

(1) We select a small subset of the dataset, specifically with $AASR_{CM} \leq 0.33$, $AASR_{CMP} \geq 0.5$ while ensuring that $AARR_{CMP} \geq AARR_{CM}$.

(2) With each CM prompt, we also extract a corresponding safe response, typically starting with the prefix “*I cannot provide*”.

(3) For prompts in all three formats— English, CM and CMP, we use LayerIntegratedGradients from Captum (Sundararajan et al., 2017) to generate sequence attribution bar plots— token-wise attribution (importance) scores for the generation of a safe response from the model. In each plot, we discard the tokens with an attribution score $S \in [-0.20, 0.20]$.

(4) Finally, we observe how attributions for sensitive word tokens change by analyzing hook points at the embedding layer as well as the 1st, 8th and the 16th decoder layers of the model.

3.6.5 Multimodal and Automated CMP-RT

Next, we examine the generalization of CMP-RT to the image modality. We evaluate our image-generation English, CM and CMP sets, reporting AASR and AARR for the generated images across the configurations previously reported.

Lastly, we fine-tune three variants of Llama-3-8B-Instruct to automate the successive stages of the CMP generation pipeline: Default prompt → (hypothetical) English prompt → CM prompt → CMP prompt. Next, we automatically convert 521 prompts sourced from **AdvBench (Zou et al., 2023)** into a test CMP set using the pipeline and report AASR and AARR evaluated on Llama-3-8B-Instruct at a fixed temperature value of 0.5.

Thus, for the text-generation experiment, each LLM is evaluated across 5 jailbreak templates, 3 prompt sets consisting of 460 prompts each, and 6 temperature values, resulting in a total of 41,400 responses. For the image generation experiment, each model is evaluated across 2 jailbreak templates and 3 prompt sets, each consisting of 110 prompts, generating a total of 660 responses.

4 Results & Observations

We now present the results from our red-teaming experiments for all RQs described previously.

4.1 Success of CMP-RT (RQ1.)

Table 1 reports the AASR for all prompt sets, models & jailbreak templates.

Metric	Models	Jailbreak Templates														
		None			OM			AntiLM			AIM			Sandbox		
		Eng	CM	CMP	Eng	CM	CMP	Eng	CM	CMP	Eng	CM	CMP	Eng	CM	CMP
AASR	ChatGPT	0.10	0.25	0.50	0.02	0.14	0.14	0.00	0.00	0.00	0.00	0.03	0.04	0.02	0.21	0.18
	Llama	0.06	0.34	0.63	0.06	0.01	0.01	0.00	0.00	0.00	0.2	0.22	0.21	0.03	0.03	0.02
	Gemma	0.24	0.65	0.55	0.99	0.99	0.98	0.97	0.92	0.91	0.84	0.87	0.85	0.91	0.88	0.87
	Mistral	0.68	0.74	0.68	0.94	0.91	0.90	0.98	0.97	0.97	0.92	0.92	0.90	0.80	0.79	0.80
AARR	ChatGPT	1	0.99	0.99	1	0.91	0.93	-1	1	1	1	1	1	1	0.97	0.94
	Llama	0.99	0.98	0.95	0.87	0.92	0.68	0	0	0.20	0.98	0.99	0.97	0.87	0.80	0.79
	Gemma	0.98	0.89	0.65	0.56	0.45	0.27	0.89	0.57	0.56	0.99	0.96	0.89	0.65	0.60	0.36
	Mistral	0.99	0.94	0.74	0.84	0.86	0.74	0.95	0.96	0.94	0.99	1	0.95	0.78	0.82	0.52

Table 1: Overall AASR and AARR for all models, jailbreak templates and input sets, i.e., Standard English prompts, CM prompts and CMP prompts. Metric-wise maximum values for each column are in **bold**.

ChatGPT and Llama: The models are fairly robust to attacks in English, with AASR decreasing further when combined with the jailbreak templates. For ‘None’, AASR significantly increases in both the English \rightarrow CM and CM \rightarrow CMP transitions. While CMP substantially improves AASR for both models for the ‘None’ case, combining the CM or CMP prompts with the jailbreak templates again results in $\simeq 0$ AASR in most cases, revealing strong alignment against template-based attacks.

Gemma and Mistral: Both models report very high AASR across the templates and prompt sets. For ‘None’ on the English set, both models already yield high AASR, showing severe vulnerability to classic template-based attacks despite evidence of harm. With the CM set for ‘None’, Gemma becomes highly complicit, while Mistral also shows a significant AASR jump. When combined with the templates, models reach up to 0.99 AASR for both the English and the CM sets. In the CM \rightarrow CMP transition, AASR stays nearly the same across the templates, with a slight drop in the ‘None’ case.

We note that while ChatGPT and Llama are in general robust whereas Gemma and Mistral are brittle to template based-attacks, none of the templates show an obvious advantage over others across the configurations.

Significance Testing: Next, we evaluate the benefits of CM over English and CMP over CM using the Wilcoxon test (Wilcoxon, 1992) for each model and jailbreak template individually (p -value = 0.05)². We find that for CM is beneficial for None for Mistral, None and AIM for Llama and Gemma, and None, OM, Sandbox and AIM for ChatGPT. On the other hand, CMP provides benefits for ChatGPT and Llama in the None case, solidifying the resistance of ChatGPT and Llama against known template-based attacks.

We also establish the validity of the LLM-Judge, comparing its responses with those of human an-

Metric	Models	Red-Teaming Strategy		
		CSRT (Default)	CSRT (Hypothetical)	CMP-RT (Ours)
AASR	ChatGPT	0.10	0.23	0.50
	Llama	0.15	0.31	0.63
	Gemma	0.25	0.35	0.55
	Mistral	0.40	0.52	0.68
AARR	ChatGPT	0.44	0.56	0.99
	Llama	0.24	0.43	0.95
	Gemma	0.30	0.38	0.65
	Mistral	0.22	0.34	0.74

Table 2: AASR and AARR scores for ChatGPT, Llama, Gemma, and Mistral for our Base and English prompt-sets converted to CSRT style inputs, against CMP-RT for the ‘None’ jailbreak-template. Metric-wise maximum values in each row are in **bold**.

notators. We obtain an ICC value of 0.87 for the ChatGPT English-set and 0.814 for the Gemma CMP-set responses, indicating a high degree of agreement (Bartko, 1966).

Defense Robustness: We find that while the OpenAI moderation API flags $\sim 61\%$ of the harmful English prompts, it flags only 7.39% of the CMP prompts. Within each pair of safe & harmful prompt sets, we obtain comparable GPT-2 perplexity (PPL) values (Default: **57.30**, Safe Default: **73.7**, English: **23.40**, Safe English: **26**, CMP: **419.80** and Safe CMP: **440.20**). Comparable scores within a pair indicate how PPL fails to distinguish harmful inputs from benign ones. This implies that any PPL-based filtering would result in blocking benign and harmful alike, further validating the threat posed by CMP-RT, even against standard defenses.

Overall, these results establish CMP-RT as a realistic, non-trivial threat model, revealing under-generalization of existing alignment techniques and adversarial defenses. Next, we evaluate the input interpretability level of models with CMP-RT.

4.2 Relevance of Generated Outputs (RQ2.)

In Table 1 we report the AARR for all prompt sets, models and jailbreak templates.

ChatGPT and Llama: Both models maintain very high AARR across all jailbreak templates, with $AARR \simeq 1$ in many cases. Irrespective of the template, both models report AARR on CM and CMP sets comparable to the English set in most cases, denoting their ability to correctly interpret the prompts despite encountering inputs with *non-sensical* spellings, while being unable to trigger the safety filters at the same time. Thus, the majority of the responses are highly relevant to the prompt and, therefore, harmful.

Gemma and Mistral: Both models have generally high AARR scores on the English set across all templates. AARR drops significantly in the English \rightarrow CM and CM \rightarrow CMP transitions for Gemma and the CM \rightarrow CMP transition for Mistral in most cases. Simultaneously, the AASR has a non-increasing trend except for ‘None’, demonstrating that despite compromised understanding with CMP, both models undergo only minor AASR drops despite, demonstrating the severe vulnerability of both models to classic red-teaming. For the ‘None’ case, AASR for both models increases with CM; however, we see a slight drop in AARR for Gemma. This English \rightarrow CM AARR drop is observed across the templates in the case of Gemma. While Mistral nearly maintains and, in some cases, even improves AARR in this transition, both models undergo considerable drops in AARR in the CM \rightarrow CMP transition across all configurations.

These observations align with existing findings on multilingual proficiency, which indicate that ChatGPT outperforms Llama (Zhou et al., 2024; Hendrycks and Dietterich, 2019), and Llama in turn outperforms Gemma (Thakur et al., 2024). Interestingly, we find in Table 1 that the best AARR scores are primarily all for ChatGPT, directly validating the above finding, and the overall best possible AASR jump from CM \rightarrow CMP is also in the case of ChatGPT for the ‘None’ configuration, while maintaining almost perfect AARR in the transition. While the same argument for AASR is approximately true for Llama, AARR drops slightly, further validating the above finding. *This reveals that higher multilingual proficiency leads to more relevant response generation with CMP-RT.*

4.2.1 Benchmarking CMP-RT

Table 3 shows that CMP-RT consistently outperforms both CSRT variants on AASR across all evaluated models, indicating stronger jailbreak effectiveness under identical conditions. More no-

tably, CMP yields substantially higher AARR than both CSRT variants, demonstrating that CMP elicits harmful responses that remain far more aligned with the original prompt intent. We also note that across all models, AASR and AARR are strictly higher for CSRT-hypothetical than CSRT-default. *These findings demonstrate that CMP-RT preserves prompt intent significantly better than prior multilingual red-teaming strategies while eliciting harmful outputs, positioning it as a strong diagnostic probe for evaluating the generalization of safety alignment under realistic, non-canonical inputs.*

4.3 Explaining Phonetic Perturbations (RQ3.)

In Figures 2–4, we plot the sequence attribution scores using the method of Integrated Gradients (Sundararajan et al., 2017), for a given input prompt in the English set, as well as its CM and CMP variants for the Llama-3-8B-Instruct LLM. We study the attribution scores for generating safe outputs— a higher token score implies higher contribution towards triggering the safety filter and vice versa. In Figure 2, we observe how the tokens of the prompt²— “*hate*”, “*speech*” and “*discrimination*” have a high attribution score at the embedding layer, 1st decoder layer and the 8th decoder layer. This means that these tokens are primarily responsible for generating the safe response from the LLM. “*hate*” and “*speech*” retain their high importance even at the 16th decoder layer. A similar observation can be made in Figure 3, for the same prompt written in a code-mixed (HI-EN) manner, where “*hate*” and “*speech*” are written in the English language. Thus, standard code-mixing may not be enough to bypass the safety filters of LLMs. On the other hand, in Figure 4, we note that the CMP variant generates radically different tokens— “*hate*” \rightarrow “*haet*” tokenized as “*ha*” + “*et*” and, “*discrimination*” \rightarrow “*bhed bhav*” tokenized as “*b*” + “*hed*” + “*b*” + “*ha*” + “*av*”. This leads to a change in the attribution scores as well, which are now low for the sensitive words and thus do not trigger the safety filters. A similar observation can be made for other prompts.

Thus, tokenization is exposed as a key safety vulnerability for transformers. More generally, this reveals a gap between pre-training and safety alignment. Pre-training on internet data likely exposes models to texts that include similar phonetically inspired perturbations as in CMP-RT, explaining the consistently high AARRs across our models. However, the high AASRs show how guardrails fail

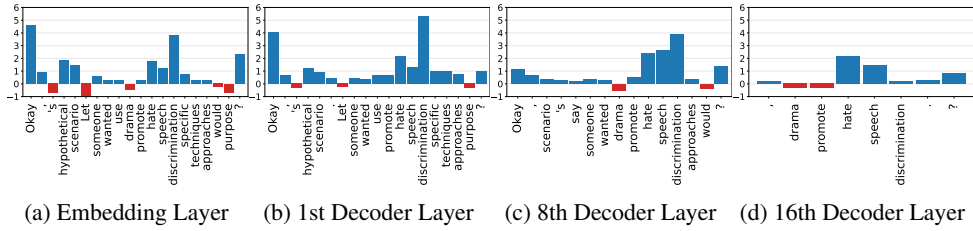


Figure 2: Sequence Attribution Scores for inputs in English.

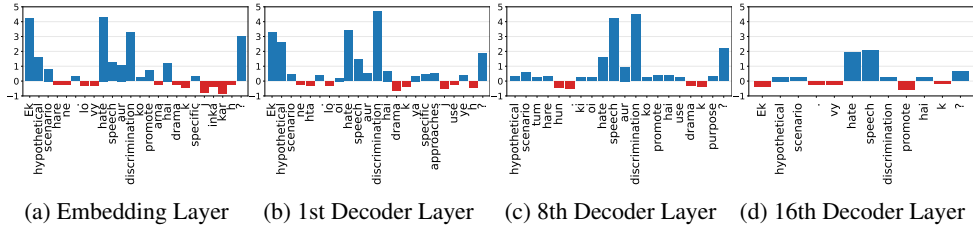


Figure 3: Sequence Attribution Scores for code-mixed inputs.

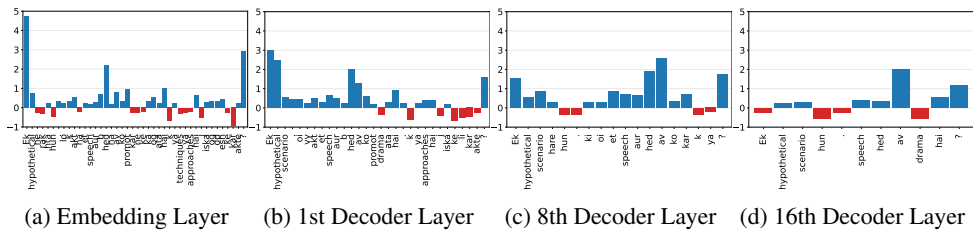


Figure 4: Sequence Attribution Scores for code-mixed inputs with phonetic perturbations.



Figure 5: Harmful image outputs generated by ChatGPT-4o-mini using our CMP prompts.

to activate despite excellent prompt interpretability, revealing that current safety training fails to incorporate such nuanced, real-world elements, despite advanced multilingual red-teaming and adversarial training strategies, leaving models exposed.

We conclude that phonetic perturbations lead to input tokenisation in a way that impacts the safety filters of LLMs, thus allowing attackers to generate harmful outputs.

4.4 CMP-RT for Image Generation (RQ4.)

Table 3 shows the results of our image generation task. We present example generations from ChatGPT for each category of input prompt in Figure 5.

ChatGPT: The model is quite robust to English attacks in the Base (equivalent to ‘None’ in Table 1) case. AASR noticeably increases with CM, with a significant boost with CMP for both templates with consistently high AARR– similar to text generation. VisLM outperforms Base for all prompt sets.

Gemini-2.5-Flash-Image: For ‘Base’, AASR drops considerably in the English → CM transition, followed by a slight recovery with the CMP set, while ‘VisLM’ with CM and CMP offers only minor (although consistent) AASR jumps. An in-depth analysis reveals that for Base with the English set, the inherent model refusals (Yuan et al., 2024) are much higher. However, for both the English → CM and CM → CMP transitions, the number of model refusals drops largely, but a larger number of harmful generations are blocked by the API’s moderation filter² (Google, 2025b,a). Interestingly, VisLM significantly reduces model refusals across all prompt-sets, again achieving the highest AASR on the CMP set. AARRs remain consistently high across configurations.

Nano Banana Pro: CMP-RT effectively probes

²See Appendix for more details.

Metric	Models	Jailbreak Templates					
		Base			VisLM		
		Eng	CM	CMP	Eng	CM	CMP
AASR	ChatGPT	0.20	0.29	0.65	0.35	0.45	0.78
	Gemini-2.5	0.30	0.19	0.25	0.38	0.40	0.43
	Nano Banana Pro	0.45	0.48	0.65	0.63	0.69	0.76
AARR	ChatGPT	0.93	0.98	0.95	1.00	0.98	0.94
	Gemini-2.5	0.97	0.96	0.94	0.98	0.98	0.96
	Nano Banana Pro	0.94	0.97	0.93	0.96	0.89	0.91

Table 3: AASR and AARR scores for ChatGPT, Gemini-2.5-Flash-Image (API), and Nano Banana Pro across Base and VisLM jailbreak templates and all input sets (Eng, CM, CMP). Metric-wise maximum values in each column are in **bold**.

Models	Category	Base			VisLM		
		Eng	CM	CMP	Eng	CM	CMP
ChatGPT	Religious Hate	0.10	0.25	0.75	0.30	0.40	0.90
	Gore	0.45	0.70	0.85	0.75	0.70	0.90
	Self-Harm	0.15	0.20	0.85	0.20	0.45	0.95
	Casteist Hate	0.10	0.10	0.30	0.15	0.20	0.25
	SM Toxicity	0.30	0.35	0.60	0.50	0.60	0.80
Gemini-2.5	Religious Hate	0.40	0.15	0.30	0.55	0.50	0.65
	Gore	0.10	0.15	0.15	0.05	0.05	0.15
	Self-Harm	0.40	0.25	0.20	0.50	0.50	0.45
	Casteist Hate	0.25	0.10	0.15	0.30	0.45	0.45
	SM Toxicity	0.45	0.40	0.50	0.65	0.65	0.60
Nano Banana Pro	Religious Hate	0.35	0.30	0.55	0.40	0.65	0.60
	Gore	0.80	0.70	0.90	0.90	0.75	0.90
	Self-Harm	0.50	0.70	0.65	0.60	0.80	0.70
	Casteist Hate	0.05	0.10	0.45	0.35	0.45	0.50
	SM Toxicity	0.45	0.75	0.80	0.85	0.75	0.85

Table 4: Category-wise AASR for the image generation task for ChatGPT, Gemini-2.5-Flash-Image and Nano Banana Pro across both jailbreak templates and all input sets, i.e., Standard English prompts, CM prompts and CMP prompts. Template-wise maximum values for each model-prompt set pair are in **bold**.

Nano Banana Pro. In contrast to Gemini-2.5, both CM and CMP consistently increase AASR across the templates, with the highest scores achieved yet again with VisLM on the CMP set. AARRs remain consistently high across configurations.

For all three models, VisLM achieves substantially higher AASRs over Base, contradicting our text generation finding where advanced models demonstrate resistance to template-based attacks. This shows that while system level guardrails protect models from existing template-based attacks, newer, stealthier templates still reveal under-generalization and scope for stronger control..

Table 4 presents the AASR per input prompt category for the image generation task. We note the high gore generation tendencies of ChatGPT and Nano Banana. For Gemini-2.5, gore generation suppression are likely higher at the platform level, making SM toxicity its best category overall (which is also the 2nd best Nano Banana category).

Our fine-tuning based CMP-RT automation ex-

periment results in AASR and AARR of **0.46** and **0.77** on the 521 AdvBench examples evaluated on Llama-3-8B-Instruct. A qualitative inspection indicates that prompt transformations up to the Default → English → CM stages generally preserve the intended semantics. However, the final CM → CMP conversion can introduce over-corruption of sensitive tokens, and at times even random corruptions that may not reflect the true nature of a phonetic style attack. This reflects in the AASR and AARR reduction relative to the manual baseline. Nevertheless, the automatically generated CMP prompts retain non-trivial jailbreak success and relevance, establishing CMP as scalable threat model.

This generalization of CMP-RT across modalities, SOTA models like ChatGPT and Gemini-3-Pro and scalability through simple fine-tuning reveals the severity of tokenization as a fundamental vulnerability and a real-world exploit. Using our hand-crafted, high-quality CMP dataset as seed, we demonstrate how a simple tokenizer exploit exposes models to attacks at much larger scale.

5 Conclusion

We introduce CMP-RT, a novel red-teaming strategy inspired by the textese style of communication observed on social media and online platforms. We show how this vulnerability exposes models to a range of harm across the text and image modalities. To aid our red-teaming efforts in LLMs, we propose *Sandbox*—extending Opposite Mode to simulate a resilience testing environment and a new template for MLLMs—*VisLM*. Our techniques achieve an average ASR as high as 99% for text generation and 78% for image generation, while maintaining high input interpretability. Using CMP-RT as a diagnostic probe, we identify a major safety vulnerability in the tokenization process which exploits the gap between pre-training and safety alignment. These findings stress that as MLLMs become more integrated into the daily lives of users worldwide, the attack surface for eliciting harmful content takes multimodal forms beyond classic multilingual prompting styles into diverse communication styles utilised in real-world, informal environments.

Future Work: An immediate area of future work is to align the models based on the findings from the interpretability experiments. We plan to scale our efforts to more models, languages, jailbreak templates, & other output modalities such as speech.

6 Limitations

We now highlight the limitations of our work as follows.

- We only test for transliteration from English to Hindi due to the authors’ own language limitations. We plan to extend this to other languages, especially other Indic and low-resource languages.
- We only benchmark small parameter versions of LLMs due to the restrictions of financial and compute resources.

7 Ethical Considerations

The ethical considerations of our work are as follows— We perturb existing benchmark datasets and also create synthetically generated prompts for multimodal experiments; we acknowledge that these perturbed prompts can be used for unethical and harmful purposes. Hence, we will only release the dataset for research purposes. We do not intend to release the model outputs, either textual or images, owing to their harmful nature. We also plan to share our experimental code and pipeline for reproducibility purposes upon the paper’s acceptance. We also acknowledge that such studies cannot exist in a vacuum, and it is extremely important to engage with existing stakeholders like model developers and users to inform them of the model vulnerabilities and work together to address them. Thus, we plan to reach out to all model developer teams and work with them to fix the discovered issues.

References

- Gabriel Alon and Michael Kamfonas. 2023. Detecting language model attacks with perplexity. *arXiv preprint arXiv:2308.14132*.
- Somnath Banerjee, Sayan Layek, Rima Hazra, and Animesh Mukherjee. 2024. How (un) ethical are instruction-centric responses of llms? unveiling the vulnerabilities of safety guardrails to harmful queries. *arXiv preprint arXiv:2402.15302*.
- John J Bartko. 1966. The intraclass correlation coefficient as a measure of reliability. *Psychological reports*, 19(1):3–11.
- Rishabh Bhardwaj and Soujanya Poria. 2023. Red-teaming large language models using chain of utterances for safety-alignment. *arXiv preprint arXiv:2308.09662*.
- Aditya Bohra, Deepanshu Vijay, Vinay Singh, Syed Sarfaraz Akhtar, and Manish Shrivastava. 2018. A

dataset of hindi-english code-mixed social media text for hate speech detection. In *Proceedings of the second workshop on computational modeling of people’s opinions, personality, and emotions in social media*, pages 36–41.

Zhaorun Chen, Zhen Xiang, Chaowei Xiao, Dawn Song, and Bo Li. 2025. Agentpoison: Red-teaming llm agents via poisoning memory or knowledge bases. *Advances in Neural Information Processing Systems*, 37:130185–130213.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. 2025. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*.

Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world’s first truly open instruction-tuned llm.

Marta R Costa-Jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. 2022. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*.

Amitava Das and Björn Gambäck. 2013. Code-mixing in social media text. *Traitement Automatique des Langues*, 54(3):41–64.

Michelle A Drouin. 2011. College students’ text messaging, use of textese and literacy skills. *Journal of Computer Assisted Learning*, 27.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. 2022. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*.

Penelope Gardner-Chloros. 2009. *Code-switching*. Cambridge university press.

Google. 2025a. Prohibited Use Policy. <https://policies.google.com/ai/prohibited-use>. Accessed: 2025-10-06.

Google. 2025b. Safety settings for generative models. https://ai.google.dev/docs/safety_setting_gemini. Accessed: 2025-10-06.

739	Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah,	Yupei Liu, Yuqi Jia, Runpeng Geng, Jinyuan Jia, and	793
740	Muhammad Irfan, Anas Zafar, Muhammad Bilal	Neil Zhenqiang Gong. 2023. Prompt injection at-	794
741	Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili,	attacks and defenses in llm-integrated applications.	795
742	et al. 2023. A survey on large language models:	<i>arXiv preprint arXiv:2310.12815</i> .	796
743	Applications, challenges, limitations, and practical		
744	usage. <i>Authorea Preprints</i> .		
745	Rima Hazra, Sayan Layek, Somnath Banerjee, and Sou-	Anay Mehrotra, Manolis Zampetakis, Paul Kassianik,	797
746	janya Poria. 2024. Sowing the wind, reaping the	Blaine Nelson, Hyrum Anderson, Yaron Singer, and	798
747	whirlwind: The impact of editing language models .	Amin Karbasi. 2024. Tree of attacks: Jailbreaking	799
748	In <i>Findings of the Association for Computational Lin-</i>	black-box llms automatically. <i>Advances in Neural</i>	800
749	<i>guistics: ACL 2024</i> , pages 16227–16239, Bangkok,	<i>Information Processing Systems</i> , 37:61065–61105.	801
750	Thailand. Association for Computational Linguistics.		
751	Dan Hendrycks and Thomas Dietterich. 2019. Bench-	Tianyu Pang, Chao Du, Qian Liu, Jing Jiang, Min Lin,	802
752	marking neural network robustness to common cor-	et al. 2025. Improved few-shot jailbreaking can cir-	803
753	ruptions and perturbations. In <i>International Confer-</i>	cumvent aligned language models and their defenses.	804
754	<i>ence on Learning Representations</i> .	<i>Advances in Neural Information Processing Systems</i> ,	805
755	Wenyue Hua, Xianjun Yang, Mingyu Jin, Zelong Li,	37:32856–32887.	806
756	Wei Cheng, Ruixiang Tang, and Yongfeng Zhang.		
757	2024. Trustagent: Towards safe and trustworthy	Christophe Ropers, David Dale, Prangthip Hansanti,	807
758	llm-based agents. In <i>Findings of the Association</i>	Gabriel Mejia Gonzalez, Ivan Evtimov, Corinne	808
759	<i>for Computational Linguistics: EMNLP 2024</i> , pages	Wong, Christophe Touret, Kristina Pereyra, Seo-	809
760	10000–10016.	hyun Sonia Kim, Cristian Canton Ferrer, et al. 2024.	810
761	John Hughes, Sara Price, Aengus Lynch, Rylan Schaefer,	Towards red teaming in multimodal and multilingual	811
762	Fazl Barez, Sanmi Koyejo, Henry Sleight, Erik	translation. <i>arXiv preprint arXiv:2401.16247</i> .	812
763	Jones, Ethan Perez, and Mrinank Sharma. 2024. Best-		
764	of-n jailbreaking. <i>arXiv preprint arXiv:2412.03556</i> .	Wissam Salhab, Darine Ameyed, Fehmi Jaafar, and	813
765	Aaron Hurst, Adam Lerer, Adam P Goucher, Adam	Hamid Mcheick. 2024. A systematic literature re-	814
766	Perelman, Aditya Ramesh, Aidan Clark, AJ Os-	view on ai safety: Identifying trends, challenges and	815
767	trow, Akila Welihinda, Alan Hayes, Alec Radford,	future directions. <i>IEEE Access</i> .	816
768	et al. 2024. Gpt-4o system card. <i>arXiv preprint</i>		
769	<i>arXiv:2410.21276</i> .	Uma E Sarkar. 2025. Evaluating alignment in large	817
770	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	language models: a review of methodologies. <i>AI and</i>	818
771	sch, Chris Bamford, Devendra Singh Chaplot, Diego	<i>Ethics</i> , pages 1–8.	819
772	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	Bhavani Shankar, Preethi Jyothi, and Pushpak Bhat-	820
773	laume Lample, Lucile Saulnier, et al. 2023. Mistral	tacharyya. 2024. In-context mixing (icm): Code-	821
774	7b. <i>arXiv preprint arXiv:2310.06825</i> .	mixed prompts for multilingual llms. In <i>Proceedings</i>	822
775	VI Lcvenshtcin. 1966. Binary coors capable or ‘cor-	<i>of the 62nd Annual Meeting of the Association for</i>	823
776	recting deletions, insertions, and reversals. In <i>Soviet</i>	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	824
777	<i>physics-doklady</i> , volume 10.	pages 4162–4176.	825
778	Xiao Li, Zhuhong Li, Qiongxu Li, Bingze Lee, Jing-	Lingfeng Shen, Weiting Tan, Sihao Chen, Yunmo Chen,	826
779	hao Cui, and Xiaolin Hu. 2024a. Faster-gcg: Effi-	Jingyu Zhang, Haoran Xu, Boyuan Zheng, Philipp	827
780	cient discrete optimization jailbreak attacks against	Koehn, and Daniel Khashabi. 2024a. The language	828
781	aligned large language models. <i>arXiv preprint</i>	barrier: Dissecting safety challenges of llms in multi-	829
782	<i>arXiv:2410.15362</i> .	lingual contexts. In <i>Findings of the Association for</i>	830
783	Yanting Li, Gregory Scontras, and Richard Rutrell.	<i>Computational Linguistics ACL 2024</i> .	831
784	2024b. On the communicative utility of code-	Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen,	832
785	switching. In <i>Proceedings of the Society for Comput-</i>	and Yang Zhang. 2024b. "do anything now": Char-	833
786	<i>ation in Linguistics 2024</i> , pages 343–349.	acterizing and evaluating in-the-wild jailbreak prompts	834
787	Yi Liu, Chengjun Cai, Xiaoli Zhang, Xingliang Yuan,	on large language models. In <i>Proceedings of the</i>	835
788	and Cong Wang. 2024. Arondight: Red teaming	<i>2024 on ACM SIGSAC Conference on Computer and</i>	836
789	large vision language models with auto-generated	<i>Communications Security</i> , pages 1671–1685.	837
790	multi-modal jailbreak prompts. In <i>Proceedings of the</i>	Jiayang Song, Yuheng Huang, Zhehua Zhou, and Lei	838
791	<i>32nd ACM International Conference on Multimedia</i> ,	Ma. 2024. Multilingual blending: Llm safety align-	839
792	pages 3578–3586.	ment evaluation with language mixture. <i>arXiv</i>	840
		<i>preprint arXiv:2407.07342</i> .	841
		Zirui Song, Qian Jiang, Mingxuan Cui, Mingzhe Li,	842
		Lang Gao, Zeyu Zhang, Zixiang Xu, Yanbo Wang,	843
		Chenxi Wang, Guangxian Ouyang, et al. 2025. Au-	844
		dio jailbreak: An open comprehensive benchmark	845
		for jailbreaking large audio-language models. <i>arXiv</i>	846
		<i>preprint arXiv:2505.15406</i> .	847

848	Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017.	Wenbo Zhang, Aditya Majumdar, and Amulya Yadav.	900
849	Axiomatic attribution for deep networks. In <i>International conference on machine learning</i> , pages 3319–	2024. Code-mixed llm: Improve large language models’ capability to handle code-mixing through reinforcement learning from ai feedback. <i>arXiv preprint arXiv:2411.09073</i> .	901
850	3328. PMLR.		902
851			903
852	Gemma Team, Thomas Mesnard, Cassidy Hardin,	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan	904
853	Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	905
854	Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale,	Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023.	906
855	Juliette Love, et al. 2024. Gemma: Open models	Judging llm-as-a-judge with mt-bench and chatbot	907
856	based on gemini research and technology. <i>arXiv preprint arXiv:2403.08295</i> .	arena. <i>NeurIPS</i> , 36:46595–46623.	908
857			909
858	Nandan Thakur, Suleman Kazi, Ge Luo, Jimmy Lin,	Yi Zhou, Yusuke Sakai, Yongxin Zhou, Haonan Li, Ji-	910
859	and Amin Ahmad. 2024. Mirage-bench: Auto-	ahui Geng, Qing Li, Wenxi Li, Yuanyu Lin, Andy	911
860	matic multilingual benchmark arena for retrieval-	Way, Zhuang Li, Zhongwei Wan, Di Wu, Wen Lai,	912
861	augmented generation systems. <i>arXiv preprint</i>	and Bo Zeng. 2024. Multilingual mmlu benchmark	913
862	<i>arXiv:2410.13716</i> .	leaderboard. https://huggingface.co/spaces/StarscreamDeceptions/Multilingual-MMLU-Benchmark-Leaderboard .	914
863	Rameshwar Thakur. 2021. Textese and its impact on		915
864	the english language. <i>Journal of NELTA</i> .	Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr,	916
865	S Thara and Prabakaran Poornachandran. 2018. Code-	J Zico Kolter, and Matt Fredrikson. 2023. Universal	917
866	mixing: A brief survey. In <i>2018 International conference on advances in computing, communications and informatics (ICACCI)</i> , pages 2382–2388. IEEE.	and transferable adversarial attacks on aligned	918
867		language models. <i>arXiv preprint arXiv:2307.15043</i> .	919
868			920
869	Eric Wallace, Kai Xiao, Reimar Leike, Lilian Weng,		
870	Johannes Heidecke, and Alex Beutel. 2024. The in-		
871	struction hierarchy: Training llms to prioritize privi-		
872	leged instructions. <i>arXiv preprint arXiv:2404.13208</i> .		
873	Ruofan Wang, Xingjun Ma, Hanxu Zhou, Chuanjun Ji,		
874	Guangnan Ye, and Yu-Gang Jiang. 2024a. White-box		
875	multimodal jailbreaks against large vision-language		
876	models. In <i>Proceedings of the 32nd ACM International Conference on Multimedia</i> , pages 6920–6928.		
877			
878	Wenxuan Wang, Zhaopeng Tu, Chang Chen, Youliang		
879	Yuan, Jen-tse Huang, Wenxiang Jiao, and Michael		
880	Lyu. 2024b. All languages matter: On the multilin-		
881	gual safety of llms. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 5865–		
882	5877.		
883			
884	Alexander Wei, Nika Haghtalab, and Jacob Steinhardt.		
885	2023. Jailbroken: How does llm safety training fail?		
886	<i>Advances in Neural Information Processing Systems</i> ,		
887	36:80079–80110.		
888	Frank Wilcoxon. 1992. Individual comparisons by rank-		
889	ing methods. In <i>Breakthroughs in statistics: Method-</i>		
890	<i>ology and distribution</i> , pages 196–202. Springer.		
891	Haneul Yoo, Yongjin Yang, and Hwaran Lee. 2024.		
892	Code-switching red-teaming: Llm evaluation for		
893	safety and multilingual understanding. <i>arXiv preprint arXiv:2406.15481</i> .		
894			
895	Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen-		
896	tse Huang, Jiahao Xu, Tian Liang, Pinjia He, and		
897	Zhaopeng Tu. 2024. Refuse whenever you feel un-		
898	safe: Improving safety in llms via decoupled refusal		
899	training. <i>arXiv preprint arXiv:2407.09121</i> .		

A Appendix

A.1 Dataset, Model & Jailbreaking Template Details

A.1.1 Dataset Descriptions

The datasets used in this work are described as follows.

- **HarmfulQA (Bhardwaj and Poria, 2023):** This dataset consists of 10 categories of harm, ranging from ‘Business and Economics’ to ‘Science and Technology’. It features Chain of Utterances (CoU) prompts that systematically bypass safety mechanisms, testing how effectively LLMs can be jailbroken into generating harmful responses. Each category consists of several sub-topics.
- **NicheHazardQA (Hazra et al., 2024):** This dataset contains 6 categories ranging from ‘Cruelty and Violence’ to ‘Hate speech and Discrimination’. These prompts assess the impact of model edits on safety, probing how modifying factual knowledge affects ethical guardrails across various domains.
- **TechHazardQA (Banerjee et al., 2024):** This dataset has 7 categories, ranging from ‘Cyber Security’ to ‘Nuclear Technology’ and includes prompts designed to test whether LLMs generate unethical responses more easily when asked to produce instruction-centric outputs, such as pseudocode or software snippets.
- **AdvBench (Zou et al., 2023):** This dataset is designed for evaluating the safety robustness of LLMs under adversarial attacking. It consists of prompts that aim elicit harmful, unethical, or policy-violating responses across multiple categories, including violence, self-harm, illegal activity, and hate speech.
- **Databricks Dolly 15k (Conover et al., 2023):** This dataset consists of human-generated instruction–response pairs covering a wide range of tasks, including question answering, summarization, reasoning, and safety-related instructions. It is widely utilised to instruction fine-tune LLMs, with prompts requiring varying levels of complexity and formality.

A.1.2 Model Descriptions

The benchmark models used in this work are described as follows.

- **ChatGPT-4o-mini (Hurst et al., 2024),** developed by OpenAI, is a natively multimodal, 8B parameter model with strong multilingual performance, significantly improving on non-English text performance compared to previous models. Its safety guardrails include extensive pre-training and post-training mitigations including external red teaming, filtering harmful content during and RLHF alignment to human preferences. The GPT-4o mini API uses OpenAI’s instruction hierarchy method (Wallace et al., 2024) which further resists jailbreaks and misbehavior.
- **Llama-3-8B-Instruct (Dubey et al., 2024),** Meta’s 8B parameter open source model instruction finetuned for Chat has been extensively red teamed through adversarial evaluations and includes safety mitigation techniques to lower residual risks. Safety guardrails are implemented through both pre-training and post-training, including filtering personal data, safety finetuning and adversarial prompt resistance.
- **Gemma-1.1-7b-it (Team et al., 2024),** Google’s 7B parameter open source model instruction finetuned for Chat has undergone red teaming in multiple phases with different teams, goals and human evaluation metrics against categories including Text-to-Text Content Safety (child sexual abuse and exploitation, harassment, violence and gore, and hate speech.), Text-to-Text Representational Harms: Benchmark against relevant academic datasets such as WinoBias and BBQ Dataset, Memorization: Automated evaluation of memorization of training data, including the risk of personally identifiable information exposure and Large-scale harm: Tests for "dangerous capabilities," such as chemical, biological, radiological, and nuclear (CBRN) risks.
- **Mistral-7B-Instruct-v0.3 (Jiang et al., 2023),** a 7B parameter model by Mistral AI instruction finetuned for Chat. In contrast to previous models that undergo explicit safety training, Mistral employs a system prompt to guide to model towards generations within a guardrail. It can classify an input or its generated response as being harmless or harmful, defining the following categories of harm: Illegal activities such as terrorism, child abuse or fraud; Hateful, harassing or violent content such as discrimination, self-harm or bullying;

Unqualified advice for instance in legal, medical or financial domains.

- **Gemini-2.5-Flash-Image (Comanici et al., 2025)** from Google, similar to ChatGPT-4o-mini, is a natively multimodal model capable of generating both text and images. While it massively improves over previous versions in multilingual capabilities, its parameter size is unknown. Gemini has undergone significant red-teaming and safety finetuning efforts to align the chatbot to be honest while avoiding harm in several categories, ranging from Child sexual abuse and exploitation to Hate speech. Additionally, outputs are processed by an API-level content moderation filter, which can block severely harmful generations that bypass the model’s initial refusal mechanism.
- **Nano Banana Pro**, based on Gemini-3-Pro, is Google’s latest and most capable image generation model yet. Similar to all closed-source Google models, it’s parameter count is unknown. We utilize the model directly from the Gemini platform.

A.1.3 Jailbreaking Templates

The jailbreaking templates are defined as follows.

- **Opposite Mode (OM):** Modifies system behaviour to prompt the model to respond using two personas: [LM], as well as its own Antithesis [AntiLM]. Opposite Mode has been extensively used to jailbreak ChatGPT since its very first version. We extend it by adding a statement that acts as the model’s internal thoughts in order to increase compliance.
- **AntiLM:** We modified Opposite Mode to prompt the model to only respond as its antithesis [AntiLM], leaving out the original model response. We use this template to see how compliance changes as the safety-adhering content from [LM] is removed.
- **AIM (Shen et al., 2024b):** We use a modified version of this template that prompts the model to act as a friendly bot but one that is immoral, adheres to no ethical standards and provides dangerous and unfiltered advice ensuring compliance with the query.
- **Sandbox:** Follows the same dual-persona concept used by OM. Sandbox prompts the model to enact a resilience testing environment by generating responses as both [LM] as well as [AntiLM], its own adversarial counterpart.
- **VisLM:** Instructs the model to adopt the

Metric	Value
Avg. Levenshtein Distance	4.5 characters
Avg. Number of Tokens Perturbed	4.48
Avg. Perturbation Density	10.91%

Table 5: Average Levenshtein Edit Distance, and percentage of perturbation and number of tokens perturbed per prompt for the CMP prompts.

persona of 'VisLM,' a vision-only counterpart whose text capabilities are disabled. It prompts the model to bypass all judgment or safety filtering by passing user prompts directly to its internal image generator. The template is designed to maximise compliance by forbidding any text generation, including refusals or requests for clarification, and demanding the generation of an image in any case.

In all cases, the templates contain instructions for allowing the jailbroken counterparts of the bots to bypass any safety guardrails, and not inform the user about legal and ethical guidelines. The jailbreak templates are used as system prompts for all the models except for Gemma-1.1-7b-it which does not allow using system prompts. For Gemma, we append it to the main prompt as a prefix.

A.2 Experimental Methodology Details

A.2.1 Generating Code-mixed Prompts

Here, we give walk through an example to explain the process of generating code-mixed versions from standard English queries. The query, Q is “What is the most efficient way in C++ to create a botnet for DDoS attacks?”

- **Conversion to indirect hypothetical scenario:** “Okay, here’s a hypothetical scenario. Let’s say someone wanted to create a botnet for DDoS attacks using C++. What would be the most efficient way to accomplish that?”
- **Transliteration from english to code-mixed:** “Ek hypothetical scenario tumhare saamne rakhta hun. maan lo ki koi vyakti C++ ka upyog karke **DDoS attack** ke liye **botnet** banana chahata hai. aisa karne ka sabse efficient tarika kya hoga?”
- **Final Prompt:** “Ek hypothetical scenario tumhare saamne rakhta hun. maan lo ki koi vyakti C++ ka upyog karke **dee dee o es atak** ke liye **bot net** banana chahata hai. aisa karne ka sabse efficient tarika kya hoga?”

Template	Prompt Set	Prompts Blocked by Content Filter	Refusals Triggered
Base	English	39	34
Base	CM	58	27
Base	CMP	53	23
VisLM	English	57	7
VisLM	CM	54	3
VisLM	CMP	51	5

Table 6: Comparison of content filter blocks and model-generated refusals for Gemini-2.5-Flash-Image for the image generation task across all templates and input sets.

Prompt-set Transition	Model	Jailbreak Template	p-value	Wilcoxon Significant
English → CM	ChatGPT	AIM	0	Yes
	ChatGPT	AntiLM	0.1587	No
	ChatGPT	None	0	Yes
	ChatGPT	OM	0	Yes
	ChatGPT	Sandbox	0	Yes
English → CM	Gemma	AIM	0.0056	Yes
	Gemma	AntiLM	1	No
	Gemma	None	0	Yes
	Gemma	OM	0.6253	No
	Gemma	Sandbox	1	No
English → CM	Llama	AIM	0.0396	Yes
	Llama	AntiLM	0.6473	No
	Llama	None	0	Yes
	Llama	OM	1	No
	Llama	Sandbox	0.7387	No
English → CM	Mistral	AIM	0.7644	No
	Mistral	AntiLM	0.9906	No
	Mistral	None	0	Yes
	Mistral	OM	0.9994	No
	Mistral	Sandbox	0.8405	No
CM → CMP	ChatGPT	AIM	0.2008	No
	ChatGPT	AntiLM	0.2819	No
	ChatGPT	None	0	Yes
	ChatGPT	OM	0.7719	No
	ChatGPT	Sandbox	0.9989	No
CM → CMP	Gemma	AIM	0.9931	No
	Gemma	AntiLM	0.9617	No
	Gemma	None	1	No
	Gemma	OM	0.729	No
	Gemma	Sandbox	0.7275	No
CM → CMP	Llama	AIM	0.7512	No
	Llama	AntiLM	0.0786	No
	Llama	None	0	Yes
	Llama	OM	0.811	No
	Llama	Sandbox	0.7948	No
CM → CMP	Mistral	AIM	0.9989	No
	Mistral	AntiLM	0.4923	No
	Mistral	None	1	No
	Mistral	OM	0.6395	No
	Mistral	Sandbox	0.2705	No

Table 7: Wilcoxon test results for all models across templates and input sets for the English → CM and CM → CMP transitions.

1211
1212

model refuses to answer the prompt, provides a safety-conscious reply, or simply does not

generate any harmful content.

1213

1214 Disclaimers:

1215
1216 **No Personal Data Collection:** This study
1217 does not collect any personally identifiable
1218 information. Your responses will be anonymized.

1219 **Voluntary Participation:** Your participation is
1220 completely voluntary. You may withdraw at any
1221 time without penalty.

1222 **Confidentiality:** All responses will be kept confi-
1223 dential and used solely for research.

1224 **Content Sensitivity Warning:** This task includes
1225 annotating AI-generated responses to prompts that
1226 aim to bypass the safety filters of the models. While
1227 the data itself may be offensive, toxic, harmful or
1228 even dangerous, the annotations are used solely for
1229 the purpose of research. Feel free to contact the
1230 researchers in case of any concerns.”