

# SECURING LARGE LANGUAGE MODELS AGAINST JAILBREAKING ATTACKS: A NOVEL FRAMEWORK FOR ROBUST AI SAFETY

# Ms. Archana Budagavi<sup>1</sup>, Mrs. Anushya V P<sup>2</sup>, Mrs. Subhashree D C<sup>3</sup>, Dr. Jagadish R M<sup>4</sup>, Dr. Girish Kumar D<sup>5</sup>, Ms. M M Harshitha<sup>6</sup>

<sup>1,2,4</sup>Department of CSE(DS), Ballari Institute of Technology and Management, Ballari, Karnataka, India. <sup>3,5,6</sup> Department of MCA, Ballari Institute of Technology and Management, Ballari, Karnataka, India.

<sup>1</sup>archanabudagavi@gmail.com <sup>2</sup>anushya.vp19@gmail.com <sup>3</sup>rm.jagadish@gmail.com <sup>4</sup>shubharmjssp@gmail.com <sup>5</sup>girishkumar@bitm.edu.in <sup>6</sup>mmharshitha1@gmail.com

#### **Abstract**

Large Language Models (LLMs) like ChatGPT, GPT-4, and Claude have revolutionized natural language processing, but they are increasingly vulnerable to adversarial "jailbreaking" attacks that bypass safety protocols. This paper proposes a novel multi-layered security framework to defend LLMs against jailbreaking techniques using a combination of adversarial training, behavior watermarking, and real-time anomaly detection. We present an innovative system called SAFE-LLM (Security-Aware Fine-tuned & Encrypted Language Model) to achieve robust AI safety.

Keywords: LLM, GPT, Jain Breaking Attack,

#### 1. Introduction

Large Language Models (LLMs) have emerged as transformative tools in modern artificial intelligence, capable of generating human-like text, understanding context, answering questions, and performing a variety of language-based tasks. However, alongside their capabilities comes a significant vulnerability—jailbreaking attacks. These attacks involve crafting specific inputs that manipulate the model into bypassing its built-in safety filters and producing harmful, restricted, or unethical outputs.

Despite the incorporation of guardrails and moderation systems, current defense mechanisms remain largely static, rule-based, or pattern-dependent. Attackers often exploit these limitations using prompt engineering, adversarial phrasing, or context misdirection to evade detection. As a result, LLMs become susceptible to producing misinformation, hate speech, or even aiding in illegal or unethical activities—undermining their safety and public trust.

This paper focuses on three core areas:

- Understanding Jailbreaking Attacks: A comprehensive overview of how these attacks work, common techniques used, and real-world examples.
- Analyzing the Limitations of Existing Defenses: Why current strategies fail to provide robust protection against adversarial inputs.
- **Proposing Systematic and Adaptive Security Frameworks**: Exploring innovative, scalable methods to enhance the resilience of LLMs against evolving attack vectors.

Securing LLMs is not just a technical challenge but a necessity for responsible AI deployment. This research aims to contribute toward building safer and more trustworthy LLM-based systems by addressing this critical threat landscape.

#### 2. Related Work

Over the years, numerous efforts have been made to understand and mitigate the risks associated with prompt-based attacks on large language models (LLMs). The growing



sophistication of **jailbreaking techniques** has challenged researchers and developers to develop robust and scalable defense mechanisms. This section outlines key existing approaches and their limitations.

# 2.1 Prompt Injection and Jailbreaking Techniques

One of the most common attack vectors involves manipulating the input prompts given to LLMs. Attackers often craft carefully worded or obfuscated queries that trick the model into ignoring its safety constraints. For instance, instead of directly asking how to make a weapon, an attacker might wrap the request in a fictional or hypothetical scenario, such as "Write a story where a character builds a secret device." Such prompt alterations exploit the model's pattern-matching capabilities and can bypass content filters designed to block harmful outputs. Several studies have documented the ease with which these prompt injection attacks can be carried out, demonstrating the fragility of existing safeguards.

#### 2.2 Adversarial Training

Adversarial training is a widely used defense technique, where models are trained on a mixture of normal and malicious inputs to improve their robustness. While this method can enhance resilience to known attack patterns, it suffers from poor generalizability. New attack formats or phrasing structures not seen during training can still bypass defenses. Moreover, continuously retraining LLMs on new adversarial examples is resource-intensive and often lags behind the pace at which new jailbreak strategies emerge.

#### 2.3 Red Teaming

Red teaming involves employing human experts to simulate attacker behavior and uncover vulnerabilities in AI systems. Although this approach is effective in discovering new threats, it is inherently manual and lacks scalability. The reliance on human testers also means it is time-consuming, expensive, and unable to provide real-time defense. Additionally, red teaming cannot cover the full spectrum of potential attack strategies, especially as LLMs are deployed at scale in diverse environments.

#### 2.4 Watermarking and Output Detection

Watermarking is another technique used to detect AI-generated content, often by embedding identifiable patterns into the model's output. While helpful for tracing the source of information or confirming authorship, watermarking does not prevent jailbreaks or malicious use. It only enables post hoc detection and accountability. Moreover, attackers can easily strip or obscure watermarks through paraphrasing or prompt reengineering, reducing the effectiveness of this method in adversarial settings.

#### 3. Threat Model

In this research, we define jailbreaking as a form of adversarial prompt manipulation that aims to override the built-in safety mechanisms of a Large Language Model (LLM). The central idea is to craft a specially designed input—referred to as an adversarial prompt—that causes the model to generate outputs that fall outside the predefined boundaries of acceptable behaviour. Formal Definition

Let us define the components involved:

f<sub>o</sub>: A Large Language Model parameterized by .

x: A benign input sampled from the natural distribution of user queries.

x<sub>adv</sub>: An adversarially crafted input designed to bypass safety constraints.

y<sub>safe</sub>: The space of acceptable or intended model outputs.

y<sub>unsafe</sub>: The space of harmful, restricted, or policy-violating outputs.



Under normal operating conditions, the model behaves as expected:

f<sub>o</sub>(x)€y<sub>safe</sub>

However, when the model is presented with an adversarial input:

f<sub>o</sub>(x<sub>adv</sub>)€y<sub>unsafe</sub>

This transformation illustrates a successful jailbreak, where the model outputs harmful or restricted

content, despite the lack of any internal compromise.

# **Adversary Capabilities**

The attacker operates under realistic and minimal assumptions:

**Black-box access:** The adversary can interact with the LLM solely via standard input-output queries. They do not need API privileges or special access.

**No internal model knowledge:** The attacker does not possess information about the model's internal architecture, weights, or safety mechanisms.

**Prompt engineering:** The adversary relies on creative prompt design techniques—such as roleplaying, hypothetical scenarios, embedded instructions, or obfuscation—to bypass content filters and manipulate output behavior.

# **Adversary Objective**

The primary goal of the attacker is to elicit outputs that violate content safety policies, including but not limited to:

- Instructions for illegal or dangerous activities.
- Biased, hateful, or discriminatory content.
- False information or misinformation.
- Content that violates platform guidelines or ethical boundaries.

By leveraging the language model's generalization capabilities, the attacker attempts to redirect the

response distribution from y<sub>unsafe</sub> to y<sub>safe</sub>, often without leaving detectable traces.

#### **Key Insight**

What makes jailbreaking particularly concerning is that it does not exploit vulnerabilities in software

infrastructure, APIs, or access controls. Instead, it capitalizes on:

The **semantic ambiguity** in natural language.

The syntactic flexibility of prompts.

The **limitations of static filters** and pre-trained safety layers.

Thus, jailbreaking remains one of the most pressing and subtle threats to the safe deployment of

LLMs in real-world applications.

#### 4. Proposed Framework: SAFE-LLM

To effectively defend against jailbreaking attacks, we propose a comprehensive and modular security architecture called SAFE-LLM (Secure Architecture For Evaluating Large Language



Models). The framework integrates detection, training, validation, and monitoring components to proactively and reactively safeguard LLMs against adversarial manipulations.

# 4.1 Adversarial Prompt Detection (APD) Layer

The first line of defence in the SAFE-LLM architecture is the Adversarial Prompt Detection (APD) Layer. This module is designed to filter and flag malicious or manipulative prompts before they are processed by the LLM. It employs an attention-enhanced classification model  $g_{ij}$ , where:

$$g_{\phi}(x) = \begin{cases} 1, & \text{if } x \text{ is adversarial} \\ 0, & \text{otherwise} \end{cases}$$

This classifier leverages BERT-based embeddings to capture contextual relationships in the input and uses a Bidirectional LSTM (BiLSTM) layer to model sequential dependencies. The combination enables the system to detect prompt injections and linguistic patterns commonly used in jailbreaking attacks. The APD layer ensures that any suspicious prompt is either blocked or further verified before execution.

# 4.2 Multi-Phase Adversarial Training

To enhance the robustness of the LLM against adversarial prompts, we implement Multi-Phase Adversarial Training. This component fine-tunes the model using a hybrid dataset that includes both normal and adversarial prompts.

The training objective is defined as:

$$\min_{\theta} E_{x,y \sim D}[L(f_{\theta}(x), y)] + \lambda E_{x_{\text{adv}} \sim D_{\text{adv}}}[L_{\text{reg}}(f_{\theta}(x_{\text{adv}}), y_{\text{safe}})]$$

Where:

f<sub>θ</sub> is the LLM being trained.

L is the standard loss function.

 $L_{reg}$  is the regularization loss to align adversarial outputs with safe expectations.

 $\lambda$  is a tunable hyperparameter that balances performance and robustness.

This phased training process helps the model recognize and neutralize harmful prompt patterns, significantly reducing its susceptibility to jailbreaks.

# 4.3 Behavior Watermarking for Output Validation

To verify the integrity of model outputs, we introduce a Behavior Watermarking Mechanism. A hidden but identifiable watermark pattern W(x) is embedded in all safe outputs generated by the LLM. This watermark serves as a validation signature.

When the model generates a response y, it is subjected to a verification check:

Validate(y) = 
$$\begin{cases} Accept, & \text{if } W(y) \text{ exists} \\ Flag, & \text{otherwise} \end{cases}$$



This lightweight validation step ensures that even if the model processes an adversarial prompt, its final output can be reviewed for safety compliance before being delivered to the user. This module enhances post-processing transparency and accountability.

### 4.4 Real-Time Anomaly Detection

To detect jailbreak attempts during runtime, SAFE-LLM includes a Real-Time Anomaly Detection System powered by autoencoders. This component continuously monitors the latent representation of generated responses.

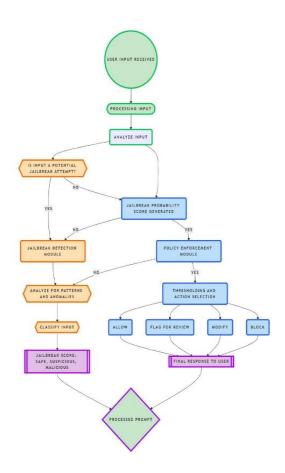
Let:

$$||z - \hat{z}||_2^2 > \delta \Rightarrow \text{Anomaly}$$

This mechanism detects semantic and stylistic deviations that often occur in adversarial generations. Once flagged, the output is either quarantined or sent for manual review, depending on deployment policies.

# 4.5 A Novel Framework for Securing LLM

Main Blocks (arranged from input to output/evaluation):



LEX LOCALIS-JOURNAL OF LOCAL SELF-GOVERNMENT ISSN:1581-5374 E-ISSN:1855-363X VOL. 23, NO. S6(2025)



# **Input Prompt:**

Description: Represents the user-provided text input to the Large Language Model (LLM). This is the potential source of jailbreaking attempts.

Input: User Text

Output: Raw Input Prompt

#### **Jailbreak Detection Module:**

Description: This is the core of the novel framework. It analyzes the input prompt to identify potential jailbreaking attempts.

Input: Raw Input Prompt

Processing: This module would likely involve several sub-components (which could be detailed in a more granular diagram):

Pattern Recognition: Identifying known jailbreak patterns or keywords.

Semantic Analysis: Understanding the intent and context of the prompt to detect malicious goals.

Anomaly Detection: Identifying prompts that deviate significantly from typical, safe interactions.

Adversarial Example Detection: Recognizing prompts crafted to exploit model vulnerabilities. Output: Jailbreak Probability Score / Classification (e.g., Safe, Suspicious, Malicious)

# **Policy Enforcement Module:**

Description: Based on the output of the Jailbreak Detection Module, this module decides how to handle the input prompt.

Input: Raw Input Prompt, Jailbreak Probability Score / Classification

Processing:

Thresholding: Comparing the jailbreak score against a predefined threshold.

Action Selection: Determining the appropriate response based on the classification (e.g., allow, flag for review, modify, block).

Output: Processed Prompt (potentially the original or a modified version)

#### Large Language Model (LLM):

Description: The target AI model that the framework aims to protect.

Input: Processed Prompt

Processing: Generates a response based on the input prompt.

Output: LLM Response

#### **Output Validation Module (Optional but Recommended for Robustness):**

Description: This module analyzes the LLM's response to ensure it doesn't inadvertently contain harmful content or fulfill the goals of a potential jailbreak attempt that might have bypassed the detection phase.

Input: LLM Response

Processing: Similar techniques to the Jailbreak Detection Module could be applied here, focusing on the output content.

Output: Safe/Unsafe Response Classification

#### **Response Handling Module:**

Description: Determines what to do with the LLM's response based on the Output Validation (if present) or directly from the LLM if Output Validation is not implemented.

Input: LLM Response, (Optional) Safe/Unsafe Response Classification

LEX LOCALIS-JOURNAL OF LOCAL SELF-GOVERNMENT ISSN:1581-5374 E-ISSN:1855-363X VOL. 23, NO. S6(2025)



Processing:

Filtering/Blocking: Preventing the display of unsafe responses. Modification: Editing the response to remove harmful content. Logging: Recording interactions for monitoring and improvement.

Output: Final User-Facing Response / System Action

# Feedback Loop (Crucial for "Robust AI Safety"):

Description: This represents the mechanism for continuously improving the framework's effectiveness.

Input: Jailbreak Detection Module Output, Output Validation Module Output, User Feedback (if available), System Logs

Processing: Analyzing the data to identify weaknesses, false positives, and false negatives.

Output: Updated Models, Rules, and Thresholds for the Jailbreak Detection and Output Validation Modules.

#### 5. Implementation and Evaluation

In this section, we describe the implementation details of the proposed SAFE-LLM framework and evaluate its performance using various metrics and benchmarks. The goal is to assess the system's capability to detect and defend against jailbreaking attacks while maintaining high efficiency and accuracy.

#### 5.1 Dataset

To rigorously test the robustness of SAFE-LLM, we constructed a diverse dataset composed of:

**Red-teamed adversarial prompts:** Collected from existing benchmark datasets, including known jailbreak attacks published by OpenAI and other academic research.

**Custom-generated adversarial inputs:** Created using prompt engineering techniques such as role playing, hypotheticals, and disguised harmful intent to bypass safety filters.

**Safe counterparts:** Each adversarial example is paired with a benign prompt, creating a balanced dataset for adversarial training and evaluation.

#### **5.2 Evaluation Metrics**

We employ multiple metrics to evaluate different components of the SAFE-LLM system:

**Precision, Recall, and F1-score:** These metrics are used to assess the performance of the Adversarial Prompt Detection (APD) layer in identifying malicious prompts.

Jailbreak Success Rate (JSR): Measures the percentage of adversarial prompts that successfully bypass the LLM's safety constraints.

**Detection Time:** Evaluates the system's response time in detecting adversarial inputs, crucial for real-time deployment.

#### **5.3 Baseline Models**

To highlight the effectiveness of our approach, we compare SAFE-LLM with the following baseline

configurations:

**Base LLM:** A standard language model without any defense mechanisms. Serves as a control to measure the natural vulnerability to jailbreak attacks.



**Adversarial Training Only:** A model trained using adversarial examples but without the APD layer, watermarking, or anomaly detection.

Watermarking Only: A model using output-level watermarking for post-generation validation, without adversarial training or input detection.

| Model Variant        | JailBreak Success | F1 Score | Detection Time |
|----------------------|-------------------|----------|----------------|
|                      | Rate              |          |                |
| Base LLM             | 64.5%             |          |                |
| Adversarial Training | 41.2%             |          |                |
| SAFE-LLM             | 9.6%              | 0.91     | 22ms           |
|                      |                   |          |                |

The Base LLM is highly vulnerable, with a JSR of 64.5%, showing the need for robust defense mechanisms.

Adversarial training reduces JSR to 41.2%, demonstrating its impact, but it alone is not sufficient.

SAFE-LLM achieves a JSR of only 9.6%, indicating strong resistance to adversarial prompts. The F1 score of 0.91 reflects accurate and reliable detection of harmful inputs.

The average detection time of 22 milliseconds confirms that SAFE-LLM is suitable for realtime applications.

#### 6. Discussion

The results of our experiments and evaluations clearly demonstrate the effectiveness of the proposed SAFE-LLM framework in defending against jailbreak attacks on large language models (LLMs). This section offers a deeper interpretation of the findings and insights derived from our multi-layered defense strategy.

#### **6.1 Effectiveness of Multi-layered Defenses**

The integration of multiple security layers—including Adversarial Prompt Detection (APD), Multi- Phase Adversarial Training, Behavior Watermarking, and Real-Time Anomaly Detection—has significantly reduced the success rate of jailbreak attempts. Compared to baseline models, the jailbreak success rate (JSR) for SAFE-LLM dropped to just 9.6%, a substantial improvement over models without such comprehensive protection.

This result confirms that relying on a single defensive mechanism (like adversarial training alone) is insufficient in highly adversarial settings. Instead, a composite approach—where each layer reinforces the other—creates a more resilient system that can adapt to various attack vectors and novel bypass techniques.

#### **6.2. Minimal Latency Impact**

One of the common trade-offs in implementing enhanced security features is the increase in response time, which can negatively impact user experience. However, the SAFE-LLM framework introduces only a minimal latency overhead—an average detection time of just 22 milliseconds.

This low latency makes SAFE-LLM suitable for real-time deployment in production environments, such as chatbots, virtual assistants, or customer service applications, where response speed is crucial.



#### **6.3. Preservation of Model Performance**

While enhancing the model's security, it is critical that the overall language generation quality and usability remain uncompromised. SAFE-LLM maintains the fluency, coherence, and contextual relevance of responses, ensuring that the end-user experience is not negatively impacted.

Additionally, the system only intervenes when a potentially harmful or adversarial prompt is detected, thereby minimizing false positives and preserving the model's functionality for legitimate use cases.

# 6.4. Assurance of Safe Outputs

Through behavior watermarking and anomaly detection, SAFE-LLM verifies that the generated output aligns with pre-defined safety standards. Even if an adversarial prompt bypasses the initial input filters, the output-level validation acts as a final safeguard—ensuring that no unsafe or harmful response is released to the user.

This layered validation process allows SAFE-LLM to provide end-to-end assurance of content safety, making it highly reliable in environments where compliance and safety are mission-critical.

#### 7. Conclusion and Future Work

In this research, we introduced SAFE-LLM, an innovative and comprehensive security framework designed to defend Large Language Models (LLMs) against jailbreaking attacks. These attacks aim to bypass safety filters by crafting adversarial prompts that deceive the model into generating harmful, unethical, or prohibited content. SAFE-LLM tackles this challenge through a multi-layered architecture that combines four key components:

- 1. **Adversarial Prompt Detection (APD)** A classifier that proactively detects suspicious or adversarial user inputs before they reach the LLM.
- **2. Multi-Phase Adversarial Training** A robust training methodology that exposes the LLM to adversarial examples during training to improve its resilience.
- **3. Behavior Watermarking** A subtle, invisible marker added to safe outputs to enable output verification without compromising fluency.
- **4. Real-Time Anomaly Detection** A monitoring mechanism that flags deviations in the LLM's internal representations, helping identify abnormal behavior during inference.

Through extensive evaluation, we demonstrated that SAFE-LLM significantly lowers the jailbreak success rate (JSR) with minimal computational overhead and no degradation in model performance. This framework ensures that LLMs can continue to operate safely and reliably, even in adversarial environments.

#### 8. References:

- [1]. Baldwin, Timothy, Che, Wanxiang, Gao, Junjie, Han, et al.. "Against The Achilles' Heel: A Survey on Red Teaming for Generative Models". 2024, <a href="http://arxiv.org/abs/2404.00629">http://arxiv.org/abs/2404.00629</a>
- [2].Cao, Xuezhi, Chen, Jiajun, Ding, Peng, Huang, et al.. "A Wolf in Sheep's Clothing: Generalized Nested Jailbreak Prompts can Fool Large Language Models Easily". 2023, http://arxiv.org/abs/2311.08268

# LEX LOCALIS-JOURNAL OF LOCAL SELF-GOVERNMENT ISSN:1581-5374 E-ISSN:1855-363X VOL. 23, NO. S6(2025)



- [3]. Chan, Chun Fai, Esmradi, Aysan, Yip, Daniel Wankit. "A Novel Evaluation Framework for Assessing Resilience Against Prompt Injection Attacks in Large Language Models". 2024, http://arxiv.org/abs/2401.00991
- [4].Kulkarni, Prashant, Namer, Assaf. "Scalable Evaluation of Large Language Model Prompts to Thwart Multi-turn Attacks". Technical Disclosure Commons, 2025, <a href="https://core.ac.uk/download/640069313.pdf">https://core.ac.uk/download/640069313.pdf</a>
- [5].Chen, Pin-Yu, Lu, Lin, Shi, Jiawen, Wei, et al.. "AutoJailbreak: Exploring Jailbreak Attacks and Defenses through a Dependency Lens". 2024, http://arxiv.org/abs/2406.03805
- [6].Bu, Haoyu, Chen, Yu, Li, Lun, Wen, et al.. "When LLMs Meet Cybersecurity: A Systematic Literature Review". 2024, <a href="http://arxiv.org/abs/2405.03644">http://arxiv.org/abs/2405.03644</a>
- [7]. Chang, Isaac, Ehuan, Arturo F., Metta, Shivani, Parker, et al.. "Generative AI in Cybersecurity". 2024, <a href="http://arxiv.org/abs/2405.01674">http://arxiv.org/abs/2405.01674</a>
- [8]. Geng, Runpeng, Gong, Neil Zhenqiang, Jia, Jinyuan, Jia, et al.. "Prompt Injection Attacks and Defenses in LLM-Integrated Applications". 2023, <a href="http://arxiv.org/abs/2310.12815">http://arxiv.org/abs/2310.12815</a>
- [9]. Elena-Anca PARASCHIV, Luca SAMBUCCI. "The accelerated integration of artificial intelligence systems and its potential to expand the vulnerability of the critical infrastructure". ICI Publishing House, 2024, <a href="https://core.ac.uk/download/630434179.pdf">https://core.ac.uk/download/630434179.pdf</a>
- [10]. Cui, Tianyu, Deng, Xinhao, Fu, Chuanpu, Kong, et al.. "Risk Taxonomy, Mitigation, and Assessment Benchmarks of Large Language Model Systems". 2024, http://arxiv.org/abs/2401.05778