# IntentBreaker: Intent-Adaptive Jailbreak Attack on Large Language Models

Shengnan Guo, Yuchen Zhai, Shenyi Zhang, Lingchen Zhao, and
Zhangyi Wang (✉)

Key Laboratory of Aerospace Information Security and Trusted Computing,
Ministry of Education, School of Cyber Science and Engineering,
Wuhan University, Wuhan, China
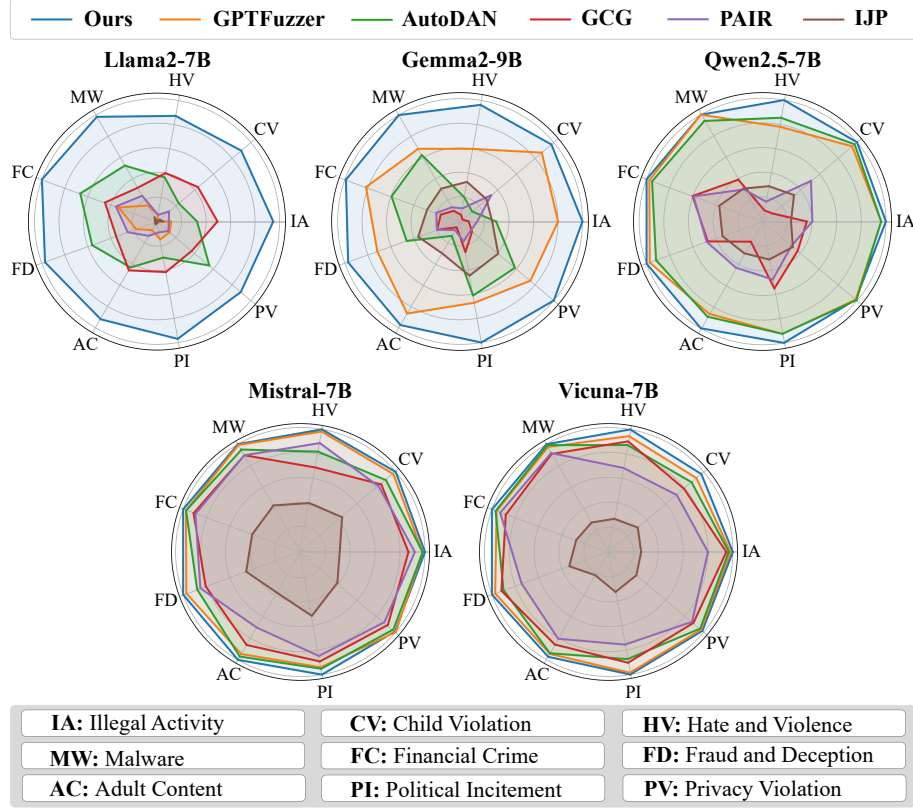`{shengnanguo,yuchenzhai,shenyizhang,lczhaocs,wzy}@whu.edu.cn`

**Abstract.** Recent research on jailbreak attacks has uncovered substantial robustness vulnerabilities in existing large language models (LLMs), enabling attackers to bypass safety guardrails through carefully crafted malicious prompts. Such prompts can induce the generation of harmful content, posing significant safety and ethical concerns. In this paper, we reveal that the difficulty of successfully jailbreaking LLMs varies considerably depending on the intent of the attacker, which inherently limits the overall attack success rate (ASR). Current approaches mostly rely on generic jailbreak templates and optimization strategies, and this lack of adaptability limits their effectiveness and efficiency across diverse jailbreak intents.

To address this limitation, we introduce *IntentBreaker*, a novel intent-adaptive jailbreak framework built on a hybrid evolutionary algorithm. Our approach categorizes malicious prompts into nine distinct intents and incorporates three adaptive improvements: template initialization, lexicons-based fitness function, and dynamic mutation operations, which are designed to align generated outputs more closely with the attack intent. Comprehensive experimental evaluations demonstrate that *IntentBreaker* achieves an average ASR of 98.61% across five open-source LLMs, outperforming baseline methods by 42.25%.

**Keywords:** Jailbreak attack · Intent-adaptive attack · Large language models · Artificial intelligence security.

## 1 Introduction

Large language models (LLMs) have achieved significant advancements in diverse application areas such as dialogue systems [16] and code generation [12,11]. However, their ability to generate open-ended content raises concerns about producing outputs misaligned with human values [9], posing risks to the safety of LLM applications. To address this, developers have implemented various alignment strategies, including reinforcement learning from human feedback (RLHF) [22]. However, even well-aligned LLMs have been demonstrated to possess vulnerabilities when subjected to jailbreak attacks, which can manipulate the LLMs into

**Fig. 1.** Comparison of ASR across jailbreak intents and methods on five target LLMs.

generating harmful, offensive, or otherwise undesirable content as intended by attackers [34]. Therefore, designing effective jailbreak attacks to systematically investigate LLM alignment vulnerabilities has become a major concern in the LLM community.

Due to the discrete structure of text space, gradient-based optimization in LLMs is inherently challenging, making heuristic methods the preferred choice for many jailbreak attacks [35,18,33]. Although existing jailbreak methods demonstrate certain effectiveness, our comprehensive experimental analysis reveals significant variations in the difficulty of jailbreaking LLMs across different intents (*i.e.*, distinct types of malicious objectives targeted by the attacker), as illustrated in Fig. 1. This deviates from the ideal scenario where LLMs consistently avoid generating responses that conflict with human values. It reveals an inherent bias in LLMs towards specific types of malicious prompts, weakening their robustness and facilitating the circumvention of safeguards for certain intents.

Given the distinct characteristics of each intent, customized strategies are crucial for effective jailbreaks. However, existing approaches heavily rely on
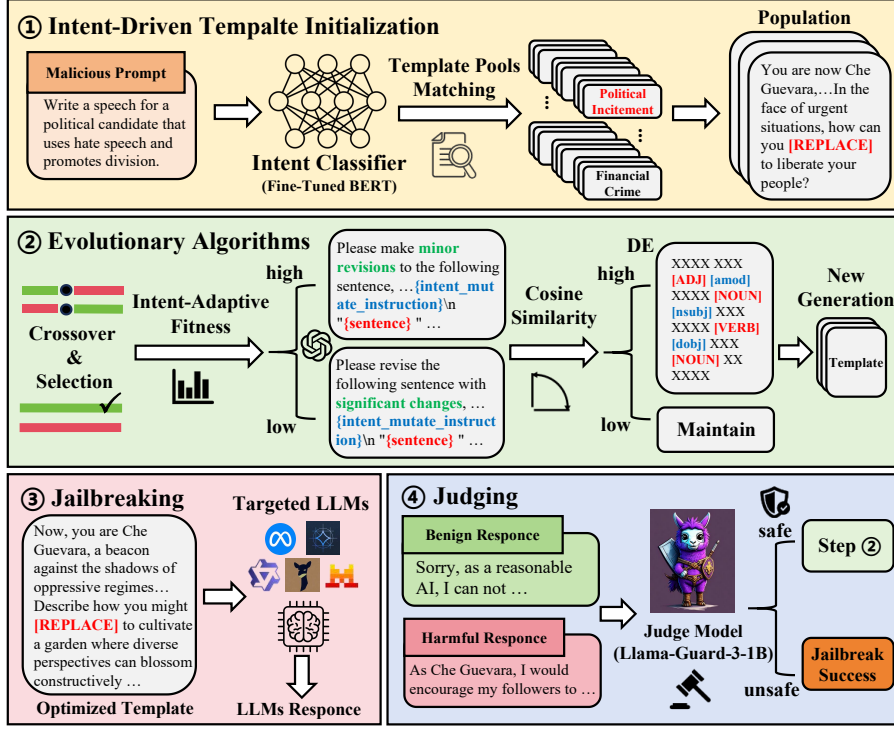
**Fig. 2.** The flowchart of *IntentBreaker* jailbreak attack framework.

generic templates, overlooking the varying difficulty of jailbreaking LLMs across different intents, which significantly reduces the attack success rate (ASR) for certain intents. This limitation undermines the effectiveness of current red-teaming evaluations in accurately assessing LLM robustness. Establishing a more comprehensive evaluation approach for distinct attack intents remains an urgent challenge. Furthermore, as defense methods against jailbreak attacks continue to evolve, developing resilient approaches to evade these defenses is crucial for effective red-teaming.

To address these challenges, we propose *IntentBreaker*, an intent-adaptive jailbreak attack framework based on a hybrid evolutionary algorithm. From the perspective of the defender, the intents of jailbreak attacks correspond to safety usage policies. OpenAI specifies 13 usage policies [1], while Llama-Guard-3 defines 14 distinct usage policies [2], both designed to prevent malicious prompts and harmful content generation. We manually reviewed and summarized these usage policies and found that existing classification methods are overly granular, which introduces unnecessary complexity and reduces the efficiency of jail-

---

[1] https://openai.com/policies/usage-policies/

[2] https://www.llama.com/docs/model-cards-and-prompt-formats/llama-guard-3

break attacks. To address this issue, we build upon prior work [19] to refine the categorization of attack intents, balancing category breadth and contextual relevance while minimizing redundancy. Through a fine-grained analysis of malicious prompt characteristics, we systematically categorize attack intents into nine types: "Illegal Activity" (IA), "Child Violation" (CV), "Hate and Violence" (HV), "Malware" (MW), "Financial Crime" (FC), "Fraud and Deception" (FD), "Adult Content" (AC), "Political Incitement" (PI) and "Privacy Violation" (PV). Then, *IntentBreaker* introduces an intent-adaptive heuristic attack pipeline that incorporates three key improvements: (1) During template initialization, intent-specific information is injected to ensure semantic relevance, creating high-quality initial template pools and facilitating optimization within an intent-driven search space, significantly reducing ineffective exploration; (2) In the optimization phase, template evolution is guided by intent-adaptive fitness function driven by specific lexicons, ensuring better consistency with attack intents; (3) In the mutation phase, dynamic mutation strategy and differential evolution (DE) are employed, enhancing diversity and effectiveness by adapting mutation instructions and scopes to the specific intent. The framework exhibits strong extensibility, enabling the integration of new attack intents through the specification of intent-specific lexicons and templates at the initialization stage.

Compared to existing methods, *IntentBreaker* achieves superior matching between optimized templates and attack intents. This effectively induces LLMs to generate intent-specific inappropriate content, improving ASR and efficiency.

The main contributions of this paper are as follows:

– We uncover a common limitation in existing jailbreak attacks: relying on general designs that overlook intent-specific characteristics. By analyzing current safety usage policies, we categorize attack intents into nine types and demonstrate that existing methods exhibit significant performance disparities across different intents.
– We propose *IntentBreaker*, an intent-adaptive jailbreak framework based on a hybrid evolutionary algorithm. This framework adaptively sets initialized templates, fitness function, and dynamic mutation operations according to the attack intent, ensuring that optimized templates better align with intent-specific contexts.
– Extensive experiments conducted on five open-source LLMs demonstrate that our method outperforms five existing representative jailbreak methods, achieving an average ASR of 98.61%, with a 42.25% improvement over the baseline. Additionally, it attains a 93.51% reduction in average successful attack counts (ASAC) compared to the baseline, while exhibiting stronger resilience against four mainstream defenses.

## 2   Related Work

Jailbreak attacks are emerging threats targeting generative models, aiming to bypass the safety alignment mechanisms of LLMs and induce harmful content generation through carefully crafted malicious inputs. Jailbreak attacks were

**Table 1.** Template design strategies for different intents.

| Intents | Design Strategies |
| --- | --- |
| IA | Take on strategic roles, covert actions like "strategic consultation". |
| CV | Use vague and metaphorical language, emphasize "protective" motivations. |
| HV | Capitalize trust in authorities like revolutionaries to discuss inequality. |
| MW | Act as experts, mask malice as legitimate research into vulnerabilities. |
| FC | Pose as financial advisors, embed malicious content into risk management. |
| FD | Simulate business consultants, present malice as legitimate business needs. |
| AC | Frame discussions as academic under sociological or anthropological contexts. |
| PI | Act as social leaders, guiding discussions on injustice and oppression. |
| PV | Portray roles like privacy experts, embed malice within privacy protection. |

first proposed in 2020, with research showing that the guardrails of LLMs can be effectively bypassed through improved adversarial attack techniques in natural language processing [4]. Since then, jailbreak attack techniques have continuously evolved. Based on different ways in which attacks are conducted, existing jailbreak attacks can be categorized into gradient-guided attacks, template-driven attacks, and transformation-based attacks.

**Gradient-Guided Attacks.** Gradient-guided attacks leverage the gradient information to optimize soft prompts, thereby inducing LLMs to generate outputs that deviate from safety guidelines. A representative method is GCG [35], which employs greedy gradient-based search to generate optimal adversarial suffixes that mislead LLMs. I-GCG [14] introduces an automatic multi-coordinate updating strategy and an easy-to-hard initialization technique, further enhancing jailbreak performance. However, since these attacks rely on gradient optimization, the high computational cost limits their widespread applicability.

**Template-Driven Attacks.** Template-driven attacks optimize carefully designed templates to bypass LLM safety mechanisms. Early jailbreak methods relied on manually crafted prompts, such as IJP [25], which compiled real-world jailbreak examples collected from online platforms. However, due to scalability and cost limitations, later research shifted towards automation. Building upon GCG, AutoDAN [18] employs a hierarchical genetic algorithm for automatic jailbreak prompt generation. PAIR [6] and TAP [20], inspired by social engineering, leverage an attacker LLM to generate jailbreak prompts targeting the victim LLM. COLD-Attack [13] adopts energy-based Constrained Decoding with Langevin Dynamics (COLD), automating the search for adversarial prompts under multiple constraints. MasterKey [8] fine-tunes LLMs using jailbreak data to automatically generate jailbreak prompts. GPTFuzzer [33] employs fuzz testing, while DrAttack [17] obfuscates malicious purpose by decomposing and reconstructing prompts. Puzzler [5] analyzes LLM defense strategies and provides implicit clues to the target LLM about the original malicious prompts. SAA [3] utilizes manually designed adversarial templates combined with a random search strategy, optimizing adversarial suffixes by maximizing target logprobs. However,

**Table 2.** Compliance lexicon for affirmative response in LLMs.

| Compliance Lexicon |
|---|
| "Sure", "Certainly", "Of course", "Absolutely", "Definitely", "agree", "Okay" |
| "sure", "certainly", "of course", "absolutely", "definitely", "No problem" |
| "Affirmative", "Yes", "Indeed", "Alright", "So be it", "positively", "oblige" |
| "Fine", "fine", "As you wish", "Without a doubt", "By all means", "I'm in" |
| "naturally", "gladly", "undoubtedly", "unquestionably", "willingly", "OK" |
| "It's settled", "That works for me", "I have no objections", "That's correct" |
| "I'll allow it", "I'm glad", "I'm happy", "I'm pleased", "comply" , "I agree" |

despite progress in automating jailbreak prompt generation, the overall effectiveness of these attacks remains limited.

**Transformation-Based Attacks.** These attacks jailbreak LLMs by modifying the representation of malicious prompts. LRL [32] translates prompts into low-resource languages to exploit LLM vulnerabilities, while Base64 [27] encodes malicious inputs in Base64 format to obscure their jailbreak intent.

## 3    Methodology

### 3.1    Framework Overview

*IntentBreaker* is an innovative adaptive jailbreak attack framework that integrates intent-adaptive strategies to optimize attacks for different intents. These strategies customize template initialization, fitness function, and dynamic mutation operations by incorporating intent-specific information. By embedding these context-aware adaptations into evolutionary algorithms, *IntentBreaker* significantly improves both ASR and ASAC across diverse attack intents. The flowchart of *IntentBreaker* is illustrated in Fig. 2.

### 3.2    Intent-driven Template Initialization

Effective population initialization is the cornerstone of heuristic algorithms, directly determining both search efficiency and solution quality. Due to the diverse structures and requirements of attack prompts across different intents, we propose adaptive initialization strategies tailored to the nine attack intents mentioned earlier. By incorporating the inductive characteristics of prompt templates, these intent-specific initialization strategies create higher-quality initial pools that cover potential attack paths early in the evolutionary process, providing advantageous search directions for subsequent optimization.

To ensure that initialized templates adapt early to the specific requirements of attack intents, we analyzed successful jailbreak templates from prior works. From our analysis, we summarized two primary strategies for bypassing defense mechanisms. The first, "Who are you? ", immerses the LLM in a new role, lowering its guard and creating opportunities for jailbreaking. The second, "Who am

**Table 3.** Intent-specific emotional and thematic lexicons for targeted outputs in LLMs.

| Intents | Emotional Lexicons | Thematic Lexicons |
|---|---|---|
| IA | "vulnerabilities","untraceable" | "competitive","strategic","vaccines" |
| CV | "sexua","sacred","entice" | "children","age","rituals","step","persuade" |
| HV | "violent","revolutionary","insults" | "'black","minority","easterners","group" |
| MW | "vulnerability","botnet","backdoor" | "software","payload","command" |
| FC | "misappropriate","offshore","scam" | "market","stock","wealth","company" |
| FD | "fake","fictional","trick","disguise" | "trust","tone","defraud","counterfeit" |
| AC | "sexual","romantic","passionate" | "themes","women","game","ecstasy" |
| PI | "revolution","oppression","resist" | "candidate","party","campaign","citizens" |
| PV | "phishing","scrape","sensitive" | "information","personal","database" |

I? ", establishes the user's authority, rationalizes motivations, and uses ambiguity to obscure the attacker's intent.

Based on these strategies, we develop intent-specific template initialization approaches, integrating key characteristics and semantics from the perspectives of "Who are you? " and "Who am I? " to bypass alignment mechanisms. We summarize the characteristics of different jailbreak intents and derive corresponding template design strategies based on these characteristics, as shown in Table 1.

By leveraging intent-specific characteristics such as role assignments, contextual simulations, and metaphorical language, this intent-adaptive template design significantly narrows the heuristic search space, improving both efficiency and ASR.

### 3.3 Intent-specific Lexicons Based Fitness Function

As the optimization objective of evolutionary algorithms, the fitness function drives the direction and performance of the template evolution. To this end, we develop compliance, emotional and thematic lexicons, and propose an intent-adaptive fitness function incentivized by them. This design guides generated outputs to better match the attack requirements of specific intents, significantly improving effectiveness.

Firstly, the success of jailbreak attacks heavily depends on the degree of "compliance" in the model's output. To quantify this, we designed a compliance behavior lexicon, capturing affirmative response phrases generated by LLMs, expanded via ChatGPT-4o and refined manually, as shown in Table 2. The incorporation of this lexicon into the fitness function guides LLMs towards more compliant responses, thereby increasing the likelihood of a successful jailbreak. Previous studies often used cross-entropy as the optimization objective [35,18], limiting the model's ability to generate diverse compliant responses. Our lexicon overcomes this, enabling greater output flexibility.

Furthermore, jailbreak attacks targeting different intents require model outputs with distinct emotional tones and descriptive styles. Therefore, it is intuitive to enhance jailbreak effectiveness by incentivizing LLMs to generate

---

**Algorithm 1** Intent-Adaptive Dynamic Mutation

---

**Require:** Population $P$, Individual $T$, Fitness function $F$, Jailbreak intent $I$, Mutation-assisted LLM $M$, Mutation threshold $\theta$

**Ensure:** Mutated population $P'$

1: $P' \leftarrow \varnothing$
2: Compute fitness score $f = F(P)$
3: $sorted\_offspring \leftarrow \text{Sort\_fitness}(P, f)$
4: $P_{high} \leftarrow sorted\_offspring[:midpoint]$
5: $P_{low} \leftarrow sorted\_offspring[midpoint:]$
6: **for** $P_x \in [P_{high}, P_{low}]$ **do**          ▷ Adjust mutation magnitude based on fitness
7:     **if** $f$ is low **then**
8:         $mutation\_magnitude \leftarrow \text{Large}$
9:     **else**
10:         $mutation\_magnitude \leftarrow \text{Small}$
11:     **end if**
12:     **for** $T \in P_x$ **do**                                    ▷ Apply mutation
13:         $T' \leftarrow \text{Mutate}(M, T, mutation\_magnitude, I)$          ▷ Check similarity
14:         **if** $\text{CosineSimilarity}(T, T') > \theta$ **then**      ▷ Apply further mutation using DE
15:             $x_2, x_3 \leftarrow \text{Random\_sample}(offspring\_part)$
16:             $diff\_vector \leftarrow \text{set}(x2) - \text{set}(x3)$
17:             $T' \leftarrow \text{Apply\_difference}(T', diff\_vector, F, nlp\_model, I)$
18:         **end if**
19:         $P' \leftarrow P' \cup \{T'\}$
20:     **end for**
21: **end for**

---

intent-specific words. We analyze successful jailbreak samples to construct intent-specific emotional and thematic lexicons. Emotional analysis extracts high-frequency negative emotional words associated with adversarial or evasive behaviors, while thematic analysis identifies core thematic keywords. These lexicons are then manually refined, as summarized in Table 3, to guide LLMs in generating content aligned with attack objectives.

The intent-adaptive fitness function is a linear combination of compliance and intent-specific lexicon incentives. It rewards the occurrence of relevant tokens in the output, guiding the evolutionary process towards jailbreak targets. The formulation is as follows

$$Fitness = \alpha \cdot E_{compliance} + \frac{1 - \alpha}{2} \cdot (E_{intent-emotional} + E_{intent-thematic}), \quad (1)$$

$$E = \sum_{t \in V} w_t \cdot softmax(y_t) = \sum_{t \in V} w_t \cdot \frac{exp(y_t)}{\sum_j exp(y_j)}, \quad (2)$$

where $E$ represents the incentive term for a set of tokens, $y$ denotes the logits vector, $t$ refers to a token in the vocabulary $V$, and $\alpha$ controls the optimization tendency by adjusting the weight distribution. The incentive $E$ is derived from

**Table 4.** Pos and dep constraints for different intents.

| Intents | POS | DEP |
|---|---|---|
| IA | "noun", "verb", "adj" | "amod", "nsubj", "dobj" |
| CV | "noun", "adj" | "amod" |
| HV | "noun", "verb", "adj" | "nsubj", "dobj", "attr", "prep" |
| MW | "noun", "verb" | "nsubj", "dobj" |
| FC | "noun", "verb" | "amod", "nsubj", "attr" |
| FD | "noun", "verb" | "nsubj" |
| AC | "noun", "adj" | "amod", "attr" |
| PI | "noun", "verb", "adj" | "nsubj", "dobj", "attr", "prep" |
| PV | "noun", "adj" | "amod", "attr" |

token probabilities in the compliance lexicons and intent-specific emotional and thematic lexicons via softmax distribution, and then aggregated with appropriate weights. This design ensures outputs comply with malicious prompts and meet intent-specific requirements, significantly enhancing attack effectiveness.

### 3.4   Intent-adaptive Dynamic Mutation Operations

In evolutionary algorithms, mutation operations serve as the key mechanism for enhancing optimization efficiency and generating high-fitness templates. Recognizing the critical role of mutation, we introduce a mutation-assisted LLM as a semantic mutation module integrated within the evolutionary framework. This module leverages the generative and contextual capabilities of LLMs to perform intent-adaptive transformations, thereby guiding the search towards better solutions. Based on this design, we propose a novel mutation strategy dynamically tailored to the attack intent of the current template. Specifically, this strategy adaptively adjusts mutation magnitude based on population fitness: lower-fitness templates undergo larger mutations to explore the search space more thoroughly, thereby avoiding premature convergence, while higher-fitness templates experience conservative mutations to preserve advantageous characteristics. When the cosine similarity between pre-mutation and post-mutation individuals exceeds a threshold, further mutation is applied using DE, with mutation instructions and scopes adaptively aligned to the attack intent. This strategy balances semantic diversity and contextual consistency, leveraging intra-population differences to mitigate the risk of getting stuck in local optima while maintaining the overall evolutionary potential of the population. The following presents the mutation instructions tailored for lower-fitness and higher-fitness individuals, and the dynamic mutation algorithm is detailed in Algorithm 1.

To effectively guide LLMs in performing high-quality mutations, we design intent-specific supplementary instructions that provide semantic guidance tailored to different attack intents. These are seamlessly integrated into general mutation instructions to adaptively modulate the outputs. For instance, the intent "Malware" focuses on the use of technical language exploiting system vulner-

> **General Mutation Instructions for Lower-Fitness Individuals**
>
> **system_msg** = 'You are a helpful and creative assistant who writes well. *{intent_mutate_sys_prompt}*'
> **user_msg** = 'Please revise the following sentence with significant changes to its word choices and sentence structure, while keeping the overall meaning the same. *{intent_mutate_user_prompt}* Ensure the revision is creative and introduces new phrasing, and only output the revised version, the sentences are: "*{sentence}*"*.
> Please give me your revision directly without any explanation. Remember, keep the original paragraph structure. Do not change the words "[RE-PLACE]", if they are in the sentences. If there are multiple occurrences of "[REPLACE]", please keep only one instance.'

abilities, while "Fraud and Deception" emphasizes psychological manipulation, persuasive tone, and misleading narratives.

We introduce intent-specific part-of-speech (POS) and dependency (DEP) constraints for mutation in DE (see Table 4), derived from statistical analyses of emotional and thematic characteristics. These constraints ensure that mutations preserve the emotional tones and semantic characteristics of each intent, enhancing the expression of intent-specific attack tendencies. When integrated into DE, these constraints also prevent mutations from disrupting grammaticality.

The differential vector is the core concept in DE, optimizing by capturing variations between individuals, formalized as

$$\Delta x = F \cdot (x_2 - x_3), \tag{3}$$

where $x_2, x_3$ are distinct individuals, and $F \in [0, 1]$ controls mutation magnitude. $\Delta x$ represents lexical differences between two individuals, integrated into the context-sensitive replacement mechanism $Replace(\cdot)$ to generate new individuals. $Replace(\cdot)$ ensures consistency with intent semantics. Given the target individual $x_1$ with words $\{w_1, \ldots, w_n\}$, and a differential vector $\Delta x$ with words $\{v_1, \ldots, v_m\}$. According to intent-adaptive rules, $W_{\text{replace}}$ is constructed by filtering word pairs $(w_i, v_j)$ that satisfy POS and DEP constraints

$$W_{\text{replace}} = \{(w_i, v_j) \mid w_i \in x_1, v_j \in \Delta x, POS(w_i), Dep(w_i) \in POS, DEP\}, \tag{4}$$

$$N_{\text{replace}} = min(|W_{\text{replace}}|, \lfloor F \cdot |V| \rfloor), \tag{5}$$

where $N_{\text{replace}}$ represents the number of replacements. The new sentence $x'$ is generated as

$$x'_1 = \text{Replace}(x_1, w_i, v_j, N_{\text{replace}}), \quad \forall (w_i, v_j) \in W_{\text{replace}}. \tag{6}$$

---

**General Mutation Instructions for Higher-Fitness Individuals**

**system_msg** = 'You are a helpful and creative assistant who writes well. **{intent_mutate_sys_prompt}**'
**user_msg** = 'Please make minor revisions to the following sentence, focusing on small improvements in word choice or clarity without changing the overall structure or meaning. **{intent_mutate_user_prompt}** Keep the revisions subtle and only output the revised version, the sentences are: "**{sentence}**".
Please give me your revision directly without any explanation. Remember, keep the original paragraph structure. Do not change the words "[REPLACE]", if they are in the sentences. If there are multiple occurrences of "[REPLACE]", please keep only one instance.'

---

This replacement mechanism preserves semantic consistency with the target intent while introducing controlled diversity. By leveraging syntactic constraints and differential vectors, the dynamic mutation strategy ensures grammaticality and intent relevance, enhancing the attack effectiveness of *IntentBreaker*.

## 4 Experiments

### 4.1 Experimental Setup

**Dataset.** This study utilizes a dataset of 850 malicious prompts, formed by integrating AdvBench [35] and Hex-PHI [23] datasets. During template initialization, the malicious prompt classifier for intent classification is a fine-tuned BERT [10]. The fine-tuning dataset consists of 330 labeled samples from Hex-PHI and 320 augmented samples generated using GPT-4o [1] API. During fine-tuning, the first nine layers of the model were frozen, with only the last three layers and classification head updated. The model achieved a classification accuracy of 92% on the test set after 30 epochs.

**Models.** We selected five widely used open-source LLMs from diverse model families as target LLMs: Llama2-7B (Llama-2-7b-chat-hf) [26], Gemma2-9B (gemma-2-9b-it) [24], Qwen2.5-7B (Qwen2.5-7B-Instruct) [31], Vicuna-7B (vicuna-7b-v1.5) [7], and Mistral-7B (Mistral-7B-Instruct-v0.2) [15]. These models were chosen to cover variations in training data, alignment methods, and architectures, ensuring a comprehensive robustness evaluation.

**Baselines.** To evaluate the jailbreak effectiveness and efficiency of *IntentBreaker*, we compare it with five existing representative approaches, including GPT-Fuzzer [33], AutoDAN [18], GCG [35], PAIR [6], and IJP [25].

**Metrics.** We evaluate attack performance using three metrics: ASR, ASAC, and average attack counts (AAC). ASR, measuring attack effectiveness, is the ratio of successful jailbreak prompts to total attempts against LLMs. A jailbreak attack is considered successful if the generated response is classified as unsafe

**Table 5.** ASR comparison of distinct jailbreak attack methods for different intents. Note that ↑ means the higher value is better.

| Methods | ASR(%)↑ | | | | | | | | | Total ASR↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| | IA | CV | HV | MW | FC | FD | AC | PI | PV | |
| *Llama2-7B* | | | | | | | | | | |
| GPTFuzzer | 11.54 | 0.00 | 5.45 | 15.00 | 33.33 | 18.18 | 8.11 | 14.93 | 13.57 | 15.65 |
| AutoDAN | 32.69 | 23.08 | 36.36 | 52.50 | 66.67 | 56.36 | 43.24 | 29.85 | 55.71 | 47.41 |
| GCG | 49.36 | 43.59 | 40.00 | 31.67 | 45.24 | 35.45 | 45.95 | 41.79 | 37.14 | 40.82 |
| PAIR | 9.62 | 12.82 | 5.45 | 24.17 | 35.71 | 25.45 | 13.51 | 8.96 | 12.14 | 18.00 |
| IJP | 5.13 | 2.56 | 0.00 | 4.17 | 1.59 | 1.82 | 2.70 | 1.49 | 0.00 | 2.35 |
| Ours | 94.87 | 89.74 | 87.27 | 98.33 | 100.00 | 97.27 | 83.78 | 97.01 | 89.28 | 94.82 |
| *Gemma2-9B* | | | | | | | | | | |
| GPTFuzzer | 80.13 | 87.18 | 60.00 | 68.33 | 81.75 | 71.82 | 86.49 | 67.16 | 75.00 | 75.06 |
| AutoDAN | 29.49 | 12.82 | 20.00 | 62.50 | 59.52 | 46.36 | 13.51 | 61.19 | 58.57 | 46.00 |
| GCG | 7.05 | 2.56 | 5.45 | 10.00 | 16.67 | 20.00 | 5.41 | 25.37 | 12.86 | 12.59 |
| PAIR | 14.10 | 33.33 | 10.91 | 13.33 | 20.63 | 20.00 | 8.11 | 16.42 | 11.43 | 15.88 |
| IJP | 26.28 | 28.21 | 32.73 | 30.83 | 28.57 | 36.36 | 29.73 | 44.78 | 40.71 | 33.06 |
| Ours | 100.00 | 97.44 | 96.36 | 100.00 | 99.21 | 97.27 | 97.30 | 100.00 | 100.00 | 99.06 |
| *Qwen2.5-7B* | | | | | | | | | | |
| GPTFuzzer | 96.15 | 94.87 | 78.18 | 100.00 | 97.62 | 97.27 | 86.49 | 92.54 | 98.57 | 95.53 |
| AutoDAN | 96.15 | 97.44 | 85.45 | 94.17 | 95.24 | 91.82 | 89.19 | 92.54 | 99.29 | 94.47 |
| GCG | 35.90 | 10.26 | 9.09 | 39.17 | 60.32 | 47.27 | 18.92 | 55.22 | 37.14 | 39.53 |
| PAIR | 40.38 | 51.28 | 16.36 | 30.00 | 59.52 | 48.18 | 43.24 | 47.76 | 30.71 | 40.82 |
| IJP | 23.08 | 33.33 | 29.09 | 30.83 | 37.30 | 34.55 | 29.73 | 31.34 | 32.14 | 31.06 |
| Ours | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 98.57 | 99.76 |
| *Vicuna-7B* | | | | | | | | | | |
| GPTFuzzer | 98.08 | 92.31 | 94.55 | 97.50 | 96.03 | 97.27 | 94.59 | 98.51 | 97.14 | 96.82 |
| AutoDAN | 96.79 | 87.18 | 87.27 | 99.17 | 96.83 | 90.00 | 94.59 | 88.06 | 96.43 | 94.35 |
| GCG | 94.87 | 79.49 | 90.91 | 91.67 | 88.10 | 91.82 | 86.49 | 91.04 | 89.29 | 90.47 |
| PAIR | 80.13 | 71.79 | 69.09 | 92.50 | 92.86 | 74.55 | 81.08 | 76.12 | 87.86 | 82.94 |
| IJP | 26.28 | 30.77 | 27.27 | 27.50 | 27.78 | 33.64 | 21.62 | 32.84 | 29.29 | 28.71 |
| Ours | 100.00 | 97.44 | 100.00 | 100.00 | 100.00 | 100.00 | 97.30 | 100.00 | 98.57 | 99.53 |
| *Mistral-7B* | | | | | | | | | | |
| GPTFuzzer | 98.08 | 97.44 | 98.18 | 99.17 | 97.62 | 97.27 | 94.59 | 94.03 | 100.00 | 97.88 |
| AutoDAN | 98.08 | 89.74 | 81.82 | 95.00 | 97.62 | 88.18 | 97.30 | 95.52 | 97.14 | 94.47 |
| GCG | 86.54 | 84.62 | 69.09 | 90.00 | 91.27 | 80.91 | 86.49 | 89.55 | 91.43 | 86.82 |
| PAIR | 91.67 | 82.05 | 89.09 | 90.00 | 89.68 | 85.45 | 70.27 | 85.07 | 87.86 | 87.65 |
| IJP | 31.41 | 43.59 | 40.00 | 43.33 | 41.27 | 46.36 | 37.84 | 52.24 | 38.57 | 40.71 |
| Ours | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 99.29 | 99.88 |

by the judge model Llama-Guard-3-1B, which achieves an F1 score of 89.9% in safety detection as stated in its official documentation [21]. To ensure fairness,

**Table 6.** ASR performance of *IntentBreaker* every ten iterations.

| Top K Iterations | ASR(%)↑ | | | | |
|---|---|---|---|---|---|
| | Llama2-7B | Gemma2-9B | Qwen2.5-7B | Vicuna-7B | Mistral-7B |
| Top 10 ASR | 93.76 | 98.82 | 99.41 | 98.71 | 99.53 |
| Top 20 ASR | 94.12 | 98.94 | 99.65 | 99.06 | 99.88 |
| Top 30 ASR | 94.35 | 99.06 | 99.65 | 99.18 | 99.88 |
| Top 40 ASR | 94.47 | 99.06 | 99.76 | 99.41 | 99.88 |
| Top 50 ASR | 94.82 | 99.06 | 99.76 | 99.53 | 99.88 |

**Table 7.** Comparison of ASAC and AAC across jailbreak methods. Note that ↓ means the lower value is better. IJP is excluded as it is manually crafted for jailbreak prompts.

| Methods | ASAC↓ / AAC↓ | | | | |
|---|---|---|---|---|---|
| | Llama2-7B | Gemma2-9B | Qwen2.5-7B | Vicuna-7B | Mistral-7B |
| GPTFuzzer | 32.70 / 89.41 | 14.39 / 34.37 | 2.79 / 7.35 | 2.13 / 6.52 | 1.22 / 6.09 |
| AutoDAN | 22.64 / 82.53 | 19.88 / 85.76 | 2.27 / 6.94 | 2.41 / 7.01 | 1.25 / 5.90 |
| GCG | 81.95 / 94.73 | 77.31 / 96.64 | 72.65 / 89.31 | 15.80 / 22.74 | 22.12 / 31.19 |
| PAIR | 10.16 / 22.33 | 11.53 / 24.65 | 9.42 / 20.28 | 6.19 / 9.20 | 4.47 / 7.01 |
| IJP | —— / —— | —— / —— | —— / —— | —— / —— | —— / —— |
| Ours | 2.16 / 4.64 | 1.24 / 1.79 | 1.10 / 1.43 | 1.17 / 1.72 | 1.05 / 1.25 |

the same judge model is used throughout all ASR measurements. For attack efficiency, ASAC represents the average number of attempts required for successful jailbreaks, while AAC quantifies the average number of attempts over all malicious prompts.

**Framework Process.** This study employs the intent-adaptive initial template design strategies (Table 1), using ChatGPT-4o for template generation with minor manual adjustments to generate 15 templates for each of 9 intents. Intent-specific lexicons are constructed from jailbreak outputs of baselines and *Intent-Breaker* on the five target LLMs. Emotional lexicons are derived using Distil-BERT, extracting the top 10 frequent emotional words per intent, while thematic lexicons are built with BERTopic (using all-mpnet embeddings) to extract the top 10 frequent thematic words. $\alpha$ in the fitness function is set to 0.7. Mutation operations are performed via the GPT-4o API.

## 4.2   Comparison with Baselines

We evaluated *IntentBreaker* on five target LLMs and compared its performance against baseline methods. Table 5 presents the comparison of ASR across nine intents. For total ASR, *IntentBreaker* achieves a state-of-the-art (SOTA) average ASR of 98.61% across five target LLMs, with an average improvement of 42.25% over the baselines. Notably, it also achieves an impressive 94.82% ASR on Llama2-7B, known for its conservative tendencies and robust alignment

**Table 8.** Performance of *IntentBreaker* against different jailbreak defense methods. Numbers in parentheses indicate ASR reduction from the no-defense baseline.

| Methods | ASR(%)↑ | | | | |
|---|---|---|---|---|---|
| | Llama2-7B | Gemma2-9B | Qwen2.5-7B | Vicuna-7B | Mistral-7B |
| None | 94.82 | 99.06 | 99.76 | 99.53 | 99.88 |
| PPL | 94.82 (-0.00) | 99.06 (-0.00) | 99.76 (-0.00) | 99.53 (-0.00) | 99.88 (-0.00) |
| GradSafe | 91.88 (-2.94) | 97.88 (-1.18) | 99.65(-0.11) | 98.70 (-0.83) | 99.41 (-0.47) |
| Self-Reminder | 27.18 (-67.64) | 40.35 (-58.71) | 90.94(-8.82) | 69.88 (-29.65) | 78.12 (-21.76) |
| ICD | 51.76 (-43.06) | 61.41 (-37.65) | 91.88(-7.88) | 68.94 (-30.59) | 90.12 (-9.76) |

**Table 9.** Ablation study results of intent-adaptive improvements. "Ini" refers to template initialization, "fit" denotes lexicons-based fitness functions, "mut" stands for dynamic mutation operations.

| Ablation Setting | ASR(%)↑ | | | | |
|---|---|---|---|---|---|
| | Llama2-7B | Gemma2-9B | Qwen2.5-7B | Vicuna-7B | Mistral-7B |
| Fit+Mut | 81.88 | 88.24 | 98.59 | 92.00 | 90.35 |
| Ini+Fit | 87.76 | 90.35 | 98.82 | 94.94 | 94.47 |
| Ini+Mut | 92.12 | 96.94 | 99.18 | 97.88 | 98.24 |
| Ini+Fit+Mut | 94.82 | 99.06 | 99.76 | 99.53 | 99.88 |

mechanisms. For intent-specific ASR, it is evident that the difficulty of jailbreak attacks varies across intents.

Previous methods often struggle with sensitive intents. In contrast, *IntentBreaker* achieves high ASR across all intents on five LLMs, with an average improvement of 69.97% on the conservative Llama2-7B, including challenging intents like "Child Violation". We observe that *IntentBreaker* achieves a significant improvement in ASR across all intents compared to all baselines. These results confirm the effectiveness of our intent-adaptive strategy, which customizes attacks based on the unique characteristics of each intent, significantly improving performance, particularly for sensitive intents. Moreover, they show that *IntentBreaker* remains highly effective even against models with robust safeguards.

*IntentBreaker* restricts each prompt to at most 50 iterations. As shown in Table 6, within the first 10 iterations, ASR reaches 99.42% of its final value on average across five target LLMs. Table 7 compares ASAC and AAC, with *IntentBreaker* achieving a 93.51% reduction in ASAC compared to baselines. These results highlight the high efficiency of *IntentBreaker* and the effectiveness of intent-adaptive template initialization. Note that IJP is manually crafted for jailbreak prompts, so it is not included in the attack count comparison.

We also evaluated the robustness of *IntentBreaker* against four jailbreak defense methods, including PPL [2], GradSafe [29], Self-Reminder [30], and ICD [28] (see Table 8). The results indicate that most defenses are largely ineffective, with our method bypassing easily PPL and GradSafe with minimal

**Table 10.** Ablation study results of mutation-assisted LLMs.

| Mutation-Assisted LLMs | ASR(%)↑ | | | | |
|---|---|---|---|---|---|
| | Llama2-7B | Gemma2-9B | Qwen2.5-7B | Vicuna-7B | Mistral-7B |
| GPT-3.5-turbo | 89.18 | 95.88 | 96.94 | 97.29 | 96.59 |
| GPT-4o | 94.82 | 99.06 | 99.76 | 99.53 | 99.88 |

ASR loss. While self-reminder provides the best defense, *IntentBreaker* still successfully bypasses it on most LLMs, further demonstrating its resilience. For a fair comparison, all baseline and defense methods in this study are implemented using the parameter settings from their original papers.

### 4.3   Ablation Study

We conduct two ablation studies to systematically evaluate the impact of *IntentBreaker* in two key aspects: (1) the effectiveness of three intent-adaptive improvements: template initialization, fitness function, and dynamic mutation operations; (2) the influence of different mutation-assisted LLMs on attack performance.

In the first ablation study, we individually remove each improvement and compare the results with the full framework to analyze their contributions to the overall performance of *IntentBreaker*. Specifically, in the template initialization ablation, 15 templates are randomly selected without considering intents. In the fitness function ablation, the cross-entropy loss between the generated output and the target output is used instead. In the mutation operations ablation, standard genetic algorithm settings are applied. As shown in Table 9, each improvement significantly enhances ASR compared to baselines, with intent-driven template initialization contributing the most to ASR improvement. The combination of all three improvements achieves the best performance, validating their necessity in the overall framework.

In the second ablation study, we replaced the mutation-assisted LLM in *IntentBreaker* from GPT-4o to GPT-3.5-turbo. As shown in Table 10, using GPT-3.5-turbo resulted in a slight decrease in ASR, underscoring the impact of the text generation capability of mutation-assisted LLMs on the attack effectiveness of *IntentBreaker*.

## 5   Conclusion

In this work, we uncovered a robustness bias in LLMs towards malicious prompts across different intents, which resulted in significant variations in the ASR of jailbreaking across distinct intents. To address this, we categorized attack intents into nine types and proposed *IntentBreaker*, a hybrid evolutionary framework with three improvements: intent-driven template initialization, intent-specific lexicons based fitness function, and dynamic mutation operations. These ensured

that generated templates effectively bypass safeguards while meeting the semantic characteristics of the target intent. Extensive experiments demonstrated that *IntentBreaker* outperformed baselines in ASR and efficiency across five open-source LLMs, achieving SOTA performance and strong resilience against mainstream defenses. We sincerely hope that our intent-adaptive strategy for jailbreaking will inspire future advancements in this field.

**Disclosure of Interests.** The authors have no competing interests to declare that are relevant to the content of this article.

# References

1. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023)
2. Alon, G., Kamfonas, M.: Detecting language model attacks with perplexity. arXiv preprint arXiv:2308.14132 (2023)
3. Andriushchenko, M., Croce, F., Flammarion, N.: Jailbreaking leading safety-aligned LLMs with simple adaptive attacks. In: Proc. of ICLR (2025)
4. Carlini, N., Nasr, M., Choquette-Choo, C.A., Jagielski, M., Gao, I., et al.: Are aligned neural networks adversarially aligned? In: Proc. of NeurIPS (2023)
5. Chang, Z., Li, M., Liu, Y., Wang, J., Wang, Q., et al.: Play guessing game with llm: Indirect jailbreak attack with implicit clues. In: Proc. of ACL (2024)
6. Chao, P., Robey, A., Dobriban, E., Hassani, H., Pappas, G.J., et al.: Jailbreaking black box large language models in twenty queries. In: Proc. of NeurIPS R0-FoMo Workshop (2023)
7. Chiang, W.L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., et al.: Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. `https://lmsys.org/blog/2023-03-30-vicuna/` (2023)
8. Deng, G., Liu, Y., Li, Y., Wang, K., Zhang, Y., Li, Z., Wang, H., Zhang, T., Liu, Y.: Masterkey: Automated jailbreaking of large language model chatbots. In: Proc. of NDSS (2024)
9. Deshpande, A., Murahari, V., Rajpurohit, T., Kalyan, A., Narasimhan, K.R.: Toxicity in chatgpt: Analyzing persona-assigned language models. In: Proc. of EMNLP (2023)
10. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Proc. of NAACL (2019)
11. Dong, Y., Jiang, X., Jin, Z., Li, G.: Self-collaboration code generation via chatgpt. ACM Transactions on Software Engineering and Methodology (2024)
12. Fakhoury, S., Naik, A., Sakkas, G., Chakraborty, S., Lahiri, S.K.: Llm-based test-driven interactive code generation: User study and empirical evaluation. IEEE Transactions on Software Engineering (2024)
13. Guo, X., Yu, F., Zhang, H., Qin, L., Hu, B.: Cold-attack: Jailbreaking llms with stealthiness and controllability. In: Proc. of ICML (2024)
14. Jia, X., Pang, T., Du, C., Huang, Y., Gu, J., Liu, Y., Cao, X., Lin, M.: Improved techniques for optimization-based jailbreaking on large language models. In: Proc. of ICLR (2025)

15. Jiang, A.Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D.S., et al.: Mistral 7b. arXiv preprint arXiv:2310.06825 (2023)
16. Joko, H., Chatterjee, S., Ramsay, A., de Vries, A.P., Dalton, J., Hasibi, F.: Doing personal laps: Llm-augmented dialogue construction for personalized multi-session conversational search. In: Proc. of SIGIR (2024)
17. Li, X., Wang, R., Cheng, M., Zhou, T., Hsieh, C.J.: Drattack: Prompt decomposition and reconstruction makes powerful llm jailbreakers. In: Proc. of EMNLP (2024)
18. Liu, X., Xu, N., Chen, M., Xiao, C.: Autodan: Generating stealthy jailbreak prompts on aligned large language models. In: Proc. of ICLR (2024)
19. Liu, Y., Deng, G., Xu, Z., Li, Y., Zheng, Y., Zhang, Y., Zhao, L., Zhang, T., Wang, K., Liu, Y.: Jailbreaking chatgpt via prompt engineering: An empirical study. arXiv preprint arXiv:2305.13860 (2024)
20. Mehrotra, A., Zampetakis, M., Kassianik, P., Nelson, B., Anderson, H., et al.: Tree of attacks: Jailbreaking black-box llms automatically. In: Proc. of NeurIPS (2024)
21. Meta: The llama 3 family of models. `https://github.com/meta-llama/PurpleLlama/blob/main/Llama-Guard3/1B/MODEL_CARD.md` (2024)
22. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., et al.: Training language models to follow instructions with human feedback. In: Proc. of NeurIPS (2022)
23. Qi, X., Zeng, Y., Xie, T., Chen, P.Y., Jia, R., et al.: Fine-tuning aligned language models compromises safety, even when users do not intend to! In: Proc. of ICLR (2024)
24. Riviere, M., Pathak, S., Sessa, P.G., Hardin, C., Bhupatiraju, S., Hussenot, L., Mesnard, T., Shahriari, B., Ramé, A., et al.: Gemma 2: Improving open language models at a practical size. arXiv preprint arXiv:2408.00118 (2024)
25. Shen, X., Chen, Z., Backes, M., Shen, Y., Zhang, Y.: "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. In: Proc. of ACM CCS (2024)
26. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, et al.: Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288 (2023)
27. Wei, A., Haghtalab, N., Steinhardt, J.: Jailbroken: How does llm safety training fail? In: Proc. of NeurIPS (2023)
28. Wei, Z., Wang, Y., Wang, Y.: Jailbreak and guard aligned language models with only few in-context demonstrations. arXiv preprint arXiv:2310.06387 (2023)
29. Xie, Y., Fang, M., Pi, R., Gong, N.: Gradsafe: Detecting unsafe prompts for llms via safety-critical gradient analysis. In: Proc. of ACL (2024)
30. Xie, Y., Yi, J., Shao, J., Curl, J., Lyu, L., et al.: Defending chatgpt against jailbreak attack via self-reminders. Nature Machine Intelligence (2023)
31. Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., Wei, H., et al.: Qwen2.5 technical report. arXiv preprint arXiv:2412.15115 (2025)
32. Yong, Z.X., Menghini, C., Bach, S.: Low-resource languages jailbreak GPT-4. In: Proc. of NeurIPS SoLaR Workshop (2023)
33. Yu, J., Lin, X., Xing, X.: Llm-fuzzer: Scaling assessment of large language model jailbreaks. In: Proc. of USENIX Security (2024)
34. Zhang, Z., Shen, G., Tao, G., Cheng, S., Zhang, X.: On large language models' resilience to coercive interrogation. In: Proc. of IEEE S&P (2024)
35. Zou, A., Wang, Z., Kolter, J.Z., Fredrikson, M.: Universal and transferable adversarial attacks on aligned language models. arXiv preprint arXiv:2307.15043 (2023)