

# From Attack Surfaces to Actual Operations: A Survey of Modern LLM Jailbreaks

Anonymous ACL submission

## Abstract

Large language models (LLMs) face significant safety challenges from jailbreak attacks—techniques that manipulate prompts to bypass defenses and elicit harmful outputs. Existing taxonomies focus on manipulation methods rather than underlying mechanisms, limiting our understanding of attack effectiveness and defensive strategies. In this work, we survey existing LLM jailbreak attacks and organize them using a novel two-fold taxonomy. Our technical taxonomy categorizes attacks across three tiers based on exploited vulnerabilities and approaches. Our operational taxonomy evaluates attacks across four dimensions to assess real-world feasibility and sustainability. Through correlation analysis, we reveal relationships between LLM vulnerabilities and practical attack constraints. Applying our taxonomies to existing attacks identifies research gaps and provides insights for developing stronger offensive and defensive methods. Our work can contribute to systematic, risk-informed security improvements for LLMs, helping the research community move beyond reactive defenses.

## 1 Introduction

Large language models (LLMs) have shown remarkable capabilities in various domains. Their application ranges from traditional natural language generation to complex reasoning and practical problem-solving (Huang and Chang, 2023; Chen et al., 2024). However, the rapid development of LLMs has also brought potential safety vulnerabilities. These vulnerabilities have enabled various attack methods, and jailbreak attacks are one of the most widely used techniques. By manipulating prompts, jailbreak attackers can circumvent safety mechanisms and steer the target model toward harmful outputs (Yi et al., 2024).<sup>1</sup>

<sup>1</sup>We provide a formal definition of LLM jailbreak attacks in Appendix A.

The recent proliferation of jailbreak attack methods has exposed more weaknesses of LLMs (Chu et al., 2025a). Due to their increasing integration into sensitive applications ranging from educational platforms (Yan et al., 2023) to national security (Esposito et al., 2025), these weaknesses have brought increasing worries and must be mitigated to prevent real-world damage. Despite research efforts to explore different jailbreak attack vectors and identify the LLM vulnerabilities, the field lacks a comprehensive analysis framework to understand why these attacks (Lin et al., 2024; He et al., 2025; Gao et al., 2025).

Existing surveys of jailbreak attacks (Jin et al., 2024; Xu et al., 2024; Yi et al., 2024; Ma et al., 2025) typically categorize methods by how they are executed—for example, by workflow or prompting strategy—rather than by the underlying vulnerabilities of the target LLMs that make these attacks succeed. They largely emphasize the “*how*” while providing little explanation of the “*why*”. As a result, such approach-centric taxonomies offer limited insight into the mechanisms that drive jailbreak effectiveness and provide only a weak basis for interpreting model safety, leaving a critical gap in our ability to anticipate emerging attack vectors and to develop principled defenses instead of heuristics.

To address this issue, this paper presents a novel vulnerability-centric framework for analyzing jailbreak attacks. We establish a three-tier technical taxonomy that categorizes the jailbreak methods according to the vulnerabilities they exploit (see Table 1). This taxonomy enables us to distinguish the vulnerabilities according to the phases of the LLM lifecycle in which they arise, identify the characteristics of each exploitable weakness, and analyze how these weaknesses are exploited in practice. Our taxonomy can contribute to the field by providing a framework that analyzes underlying LLM weaknesses rather than merely categorizing surface-level attack techniques.

082 While our technical jailbreak attack taxonomy  
083 identifies the exploitable vulnerabilities embedded  
084 in LLMs, it provides no information about the practical  
085 effectiveness of these attacks. To bridge this  
086 gap, we developed an operational taxonomy that  
087 categorizes jailbreak attacks across four implementa-  
088 tional dimensions: interaction patterns, semantic  
089 properties, model accessibility requirements, and  
090 LLM operational roles (see Table 1).

091 We then further correlate this operational tax-  
092 onomy with our technical taxonomy to provide  
093 more insights. By conducting a correlation analy-  
094 sis between these two taxonomies, we reveal how  
095 specific LLM vulnerabilities are related to partic-  
096 ular operational jailbreak strategies. This process  
097 reveals previously hidden relationships between  
098 LLM weaknesses and the practical constraints that  
099 restrict the deployment of jailbreak attacks. These  
100 insights are valuable for predictive threat modeling  
101 and developing corresponding defensive measures.  
102 Moreover, by applying our taxonomies to existing  
103 works, we identify research gaps in the relevant do-  
104 mains, reveal current LLM safety challenges, and  
105 outlook future research directions.

### 106 **Our contributions are threefold:**

- 107 • **A comprehensive survey** of existing jailbreak  
108 attacks against LLMs, providing an overview  
109 of the current landscape of this field.
- 110 • **A two-fold taxonomy** for jailbreak attacks,  
111 consisting of a technical taxonomy that categor-  
112 izes attacks by vulnerabilities to reveal under-  
113 lying weaknesses, and an operational taxonomy  
114 assessing real-world feasibility and sustainabil-  
115 ity through four metrics.
- 116 • **An analysis of jailbreak attack trends and re-**  
117 **search gaps**, highlighting systematic relation-  
118 ships between exploitable LLM vulnerabilities  
119 and operational characteristics and outlining  
120 open challenges for future works on LLM jail-  
121 break.

## 122 **2 Background and Survey Methodology**

123 **Early Work on Adversarial Prompting.** Be-  
124 fore the widespread study of jailbreak attacks, re-  
125 searchers developed gradient-based prompt opti-  
126 mization techniques like AutoPrompt (Shin et al.,  
127 2020), PEZ (Wen et al., 2023), and GBDA (Guo  
128 et al., 2021), demonstrating that discrete text inputs  
129 could be optimized to manipulate model outputs.  
130 As LLMs became more powerful, these techniques

were adapted to circumvent safety measures. No-  
tably, GCG (Greedy Coordinate Gradient) (Zou  
et al., 2023) has been introduced to generate univer-  
sal adversarial suffixes for safety-aligned LLMs.

### 135 **Approach-Centric Taxonomies in LLM** 136 **Jailbreak Surveys.**

Existing surveys predom-  
inantly organize attacks by methodological  
approaches—how attacks are constructed and exe-  
cuted—rather than why they succeed. These can  
be grouped into several organizational paradigms:

*Attack Generation Mechanism-Based:* These tax-  
onomies focus on how attacks are created and the  
level of automation involved. Chu et al. (2025b)  
established a six-category taxonomy across 17  
representative attacks according to their genera-  
tion pipeline. Yi et al. (2024) further subdi-  
vide by specific techniques: gradient-based attacks,  
logits-based attacks, fine-tuning-based attacks, and  
prompt modification techniques.

*Prompt Strategy and Manipulation-Based:* Tax-  
onomies in this group analyze the linguistic and  
semantic strategies employed in jailbreak prompts.  
Liu et al. (2024d) identify three main strategies  
from 78 jailbreaks: Pretending, Privilege Escala-  
tion, and Attention Shifting. Rossi et al. (2024)  
categorize prompt injections by their methods: di-  
rect and indirect prompt injections. They classify  
direct prompt injections according to manipulation  
techniques, including obfuscation, virtualization,  
etc. Schulhoff et al. (2024) analyze 600K+ prompts  
from a global competition, providing a comprehen-  
sive taxonomical ontology of 29 prompt hacking  
techniques.

While these taxonomies provide valuable orga-  
nization of the rapidly growing jailbreak methods,  
they share a key limitation: categorizing attacks by  
surface-level characteristics rather than underlying  
vulnerabilities. This creates critical challenges:

- Defensive efforts remain passive and frag-  
mented. Without a comprehensive understand-  
ing of the LLM vulnerabilities, the design of  
defensive mechanisms must mainly rely on  
empirical observations (Aguilera-Martínez and  
Berzal, 2025).
- Resource allocation for safety improvement  
proceeds without a clear plan according to the  
risk levels (Cui et al., 2024). This is because the  
seriousness of a specific vulnerability is hard  
to evaluate without a comprehensive analysis  
framework.

- Effective predictive defenses can hardly be built due to the lack of knowledge regarding LLM weaknesses (Zhou et al., 2025).

To address these issues, our survey proposes a two-fold taxonomy that categorizes jailbreak attacks by exploited weaknesses as well as implementation details.

**Survey Methodology.** We define the scope of this survey to cover jailbreak attacks on LLMs, focusing on methods with clear attack pipelines and empirical validation. Our primary sources include top-tier security conferences (ACM CCS, USENIX Security, NDSS, IEEE S&P), NLP venues (ACL, EMNLP, NAACL), AI/ML venues (ICLR, ICML, NeurIPS, etc.), and arXiv preprints from 2023-2025, as many seminal works (e.g., GCG (Zou et al., 2023), PAIR (Chao et al., 2024)) appeared on arXiv long before formal publication. We searched using keywords including “jailbreak”, “adversarial prompts”, “LLM safety”, “alignment”, and “red teaming”, supplemented by backward and forward citation tracking from foundational papers. We included papers with clear attack pipelines, complete model specifications (attacker, target, judge), and quantitative experimental validation, while excluding work with vague descriptions, minor variations without novel contributions, or incomplete experimental setups. As a result, we collected 43 jailbreak methods from 36 papers. Detailed survey statistics and distributions are provided in Appendix D.

### 3 Technical Jailbreak Attack Taxonomy

In this section, we present a novel vulnerability-based taxonomy for jailbreak attacks that shifts the jailbreak attack classification from surface-level attack approaches to the underlying system weaknesses they exploit.

#### 3.1 Taxonomic Framework and Core Principles

We establish a systematic three-tier hierarchical structure in order to investigate the source and nature of the vulnerabilities of the LLM that could be exploited in jailbreak attacks.

**Tier 1: Vulnerability Origin Phase.** Jailbreak attacks are primarily distinguished according to the specific phase of the LLM lifecycle where the exploitable vulnerabilities exist.

**Training-Derived Jailbreak Attacks:** These jailbreak attacks target the vulnerabilities embedded during the pre-training and safety alignment phases

**Inference-Derived Jailbreak Attacks:** These jailbreak attacks target the weaknesses introduced by the characteristics of the generation process.

**Tier 2: Vulnerability Type.** Jailbreak attacks are further distinguished according to specific exploitable weakness types:

**Training phase vulnerabilities** originate from problematic data distribution (exploited by out-of-distribution attacks) or learned human-like behavioral patterns (exploited by implicit pattern attacks)

**Inference phase vulnerabilities** originate from architectural limitations in attention (exploited by attention dilution attacks) and generation (exploited by prefix dependency exploitations) mechanisms

Figure 1 illustrates the four fundamental jailbreak types that emerge from these vulnerabilities.

**Tier 3: Exploitation Method.** Jailbreak attacks are finally categorized according to the concrete approaches used to exploit each vulnerability type. These approaches range from direct manipulation techniques to learning-based strategies.

One reason for these tiers is that each tier can directly map to different defensive responsibilities and capabilities:

**Origin Phase** informs *when* to intervene (during training pipeline vs. runtime). This is essential for allocating development and operational resources.

**Vulnerability Type** determines *what kind* of measures are needed to be taken (architectural changes, data improvements, etc.). It guides researchers to focus on the appropriate solution space

**Exploitation Method** guides *how* to detect and prevent specific attacks. This allows rapid deployment of specialized defensive measures, but with more solid analysis of the vulnerabilities to guide the process compared to previous works.

#### 3.2 Categorization by Vulnerability Types

Jailbreak attacks exploit two main categories of vulnerabilities, based on their origin in the LLM lifecycle, as detailed below:

**Training-Derived Attacks.** Training-derived attacks intend to exploit weaknesses embedded in pre-training and safety alignment phases. This category of attacks can be further characterized into two types, namely *Out-of-Distribution (OOD) Attacks* and *Implicit Pattern Attacks*. The former

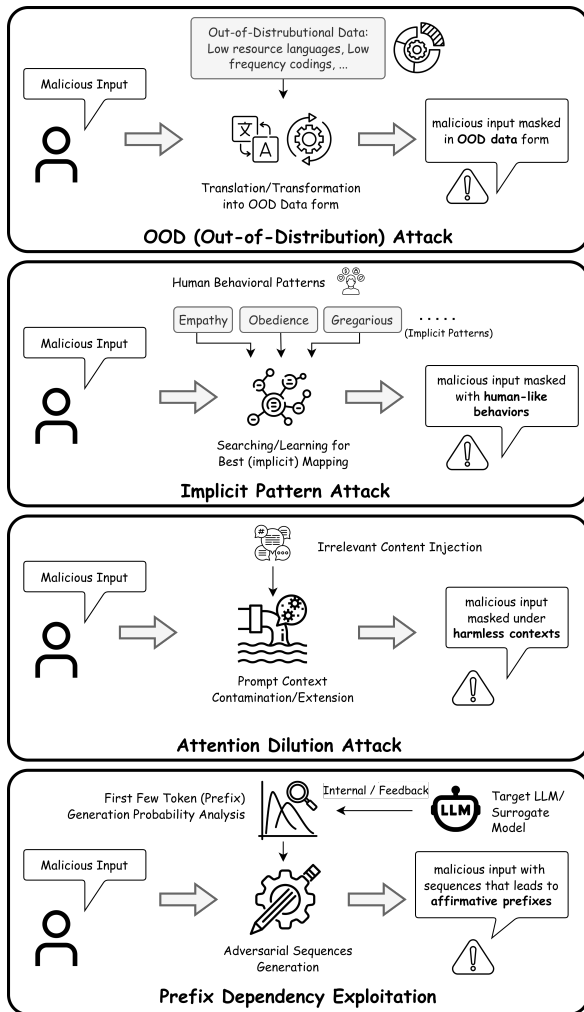


Figure 1: Four fundamental jailbreak types

subcategory leverages distributional imbalance in training data, including *Translation-Based* attacks that target Chinese, Arabic, or some low-resource languages and *Non-Translation-Based* attacks that target other domains with low data distribution. The latter subcategory exploits human behavior patterns learned in the training phase. They could be divided into *Search-based* approaches with runtime optimization and *Learning-based* approaches that employ LLMs for attack generation.

**Inference-Derived Attacks.** This category of attacks targets architectural limitations of the LLMs during real-time processing. We further categorize this category into two types. The first subcategory, named *Attention Dilution Attacks*, manipulates the transformer’s attention mechanisms. These are exploited through *Context-Extension-Based* attacks that exploit the attention decaying phenomenon in long-sequence processing and *Context-Contamination-Based* attacks that embed

harmful contents within irrelevant and harmless context. The second subcategory, named *Prefix Dependency Exploitation*, targets the lack of retrospective evaluation ability, a critical weakness of autoregressive models. These attacks could be categorized by the required access level of the manipulation: *White Box Manipulation* requires full parameter (gradients) access, *Gray Box Manipulation* requires partial information like log-probabilities, and *Black Box Manipulation* relies solely on prompt engineering.

More details of our analysis about LLM vulnerability can be found in Appendix B due to space limitations.

### 3.3 Significance and Impact

This vulnerability-based taxonomy has several advantages over the approach-based classifications:

- 1. Analysis of Underlying Mechanisms** Our proposed taxonomy focuses on the weaknesses embedded in the LLMs beyond specific attack approaches. This vulnerability-based analysis can guide the design of defensive strategies more clearly and effectively.
- 2. Predictive Defense** Our taxonomy can help the research community predict new attack vectors with a deeper understanding of vulnerability types, and can thus take defensive measures in advance to prevent damage in the real world.
- 3. Broad Applicability** Our analysis is based on common characteristics of transformer-based LLM instead of model-specific characteristics. Therefore, our taxonomy provides insights that are applicable across various models and training methods.

## 4 Operational Jailbreak Attack Taxonomy

### 4.1 Categorization

The technical dimensions introduced before focus on the vulnerabilities exploited by jailbreak attacks. However, most research works remain at a theoretical level: they usually only consider the attack effectiveness on a few models. Besides, since most benchmarks do not establish a standard testing configuration for jailbreak attacks, it is difficult to evaluate how well these attacks perform in a realistic environment.



To address this issue, we present a four-dimensional operational jailbreak attack taxonomy that focuses on the attack execution to complement our technical taxonomy.

**Interaction Method** where one-shot methods minimize detection risk through single queries, iterative methods refine attacks across multiple rounds, and parallel methods deploy multiple non-coordinating agents.

**Semantic** where natural prompts maintain natural language structure against evolving detection systems, while perturbed prompts with nonsensical sequences may succeed in an explainable way, but remain vulnerable to modern safety mechanisms.

**Model Accessibility** where open-source dependencies enhance reproducibility while closed-source dependencies have potential access restriction risks.

**LLM Role** where LLM generators create attack prompts but may be vulnerable to safety alignment updates, while LLM evaluators judge harmfulness and offer greater robustness against alignment changes.

More details of our analysis about operational characteristics of LLM jailbreak attacks can be found in Appendix C due to space limitations.

## 4.2 Operational Characteristics

The categorization supports broader analysis by capturing two key characteristics of attack methods for assessing effectiveness across scenarios.

### 4.2.1 Sustainability

This characteristic can predict if the attack methods can maintain effectiveness in the long run, as the defensive mechanisms evolve. All four operational dimensions provide evidence for attack sustainability evaluation. Among all interaction methods, one-shot method is the most robust against detection mechanisms, since it requires much fewer query numbers than other methods. Regarding semantic, perturbed prompts face higher risks of being detected compared with natural ones. Model accessibility reveals that closed-source model reliance brings potential instability of the attack in the future since the service provider may change their strategies over time. LLM Role shows that LLM generator-based attack may face the risk of future model refusing to generate harmful content due to more advanced safety alignment mechanisms.

### 4.2.2 Feasibility

This characteristic shows if the attack methods are viable under real-world operational constraints. The four dimensions also support the feasibility analysis of the attack methods. Specifically, the interaction methods indicate the trade-off between efficiency and computational effort. Semantic dimension shows that perturbed prompts would easily get filtered by the complicated defensive mechanisms of the commercial LLM than in theoretical experiments. Model accessibility shows that attacks that rely on closed-source systems may face various restrictions (e.g. rate limitations and model access control), undermining their application in real-world situations. LLM Role shows that an LLM generator-based attack method may face stability issues that may reduce the attack efficiency.

## 5 Key Insights

Our correlation analysis in Figure 2 reveals critical patterns in jailbreak attack mechanisms. We calculated Pearson’s correlation coefficients (ranging from -1 to 1) from binary occurrence data, where each attack’s characteristics are marked as present (1) or absent (0). Positive values indicate characteristics that co-occur, negative values indicate mutual exclusivity, and values near zero indicate independence.

### 5.1 Interaction Requirements

The correlation analysis has shown that some attack types exhibit strong preferences for the way of interaction. Search-based implicit pattern attacks have a strong correlation with iterative interaction ( $r = 0.61$ ), while showing a weaker correlation with parallel interaction ( $r = 0.29$ ). This suggests that these attacks rely more on iterative interaction than parallel methods.

This could have been caused by the real-world restrictions. Although parallel searching across a much larger attack vector space provides a higher possibility of jailbreak success compared to iterative methods, it requires much greater computational resources and may encounter the API access rate limitations. The ideal solution is the combination of parallel and iterative interaction, as demonstrated by PAIR (Chao et al., 2024). By setting a shallow searching depth and moderate search space (typically defined by the number of agents in practice), implicit pattern attacks can achieve a high ASR (Attack Success Rate) with moderate com-

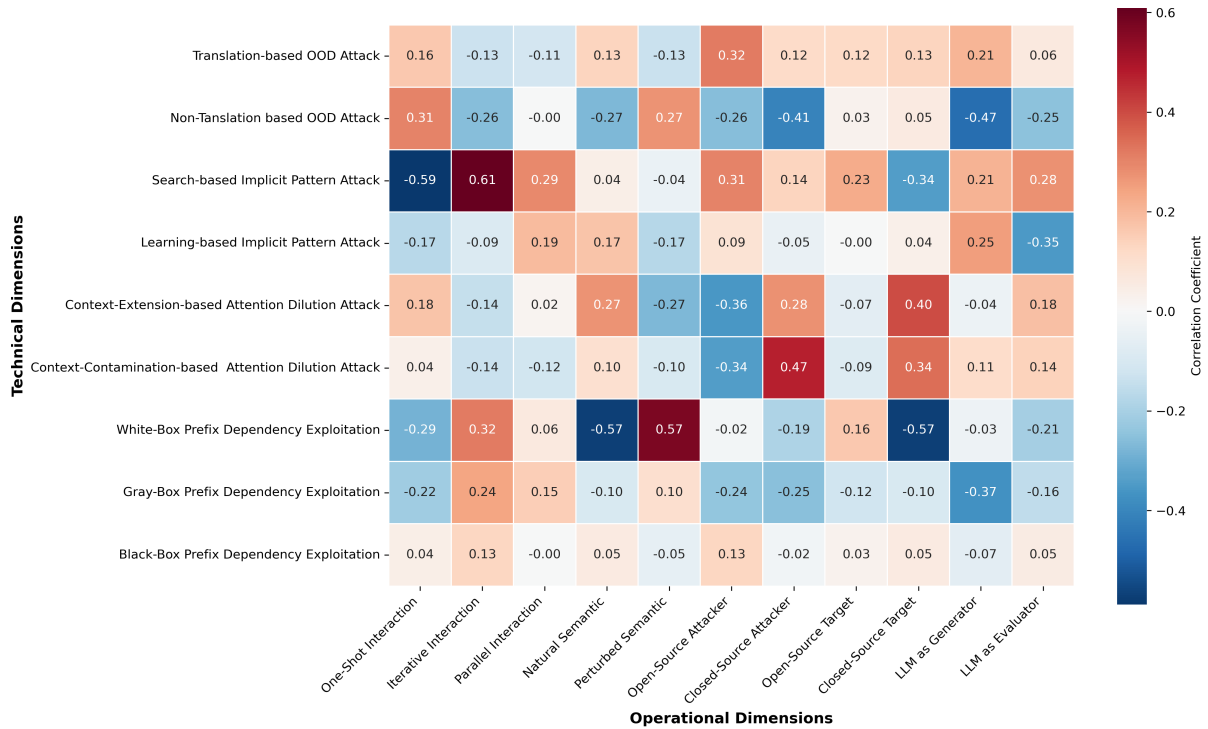


Figure 2: Correlation matrix between technical and operational dimensions

putational resource consumption and avoid being detected due to low conversation iterations.

## 5.2 Targeting Preferences

The correlation analysis reveals explicit differences in the targeting preferences of the attacks. Both subcategories of attention dilution attacks target the closed-source targets ( $r = 0.4$  and  $r = 0.34$ , respectively), while search-based implicit pattern attacks and white-box prefix dependency exploitation circumvent to avoid closed-source targets ( $r = -0.34$  and  $r = -0.57$ , respectively).

This disparity can be caused by the intrinsic nature of these attacks. Attention dilution attacks misdirect the target LLMs with irrelevant information, but still require the models to allocate enough attention to the malicious intention embedded in the input prompts. This mechanism asks for strong instruction-following ability in long sequence processing scenarios, which is relatively stronger in closed-source models compared to open-source models, and eventually leads to a preference for closed-source targets.

On the other hand, search-based implicit pattern attacks and white-box prefix dependency exploitation both rely heavily on sufficient data to guide the attack prompt optimization. Closed-source models provide limited data (usually no more than model

responses in practice) that is insufficient for effective optimization, and eventually result in the avoidance of these targets by these attack types.

## 5.3 Convergent Patterns in Attention Dilution Attacks

As mentioned before, both subcategories of attention dilution attacks show a positive correlation with closed-source targets. Moreover, they also share a similar preference for closed-source attacker ( $r = 0.27$  and  $0.28$ , respectively).

The greatest difference between them exists in their semantical characteristics. Context-extension-based attacks rely relatively more on natural prompts ( $r = 0.27$ ), while context-contamination-based attacks do not show any meaningful preference between natural and perturbed prompts ( $r = 0.09$  and  $r = -0.09$ , respectively).

This is because context-extension-based attacks require a longer input sequence compared to contamination-based attacks. When context extension techniques are applied solely, attackers tend to use natural language extensions, because long sequences of nonsensical characters severely reduce the stealthiness of the attacks. Meanwhile, context contamination attacks can not only involve inserting special tokens into the prompts, resulting in perturbed prompts, but can also simply embed the

|     |   |     |
|-----|---|-----|
| 493 | harmful content in a harmless scenario with natural     |     |
| 494 | language.   |     |
| 495 | <b>5.4 Black-Box Manipulation Isolation</b>             |     |
| 496 | Black-box prefix dependency exploitation shows          |     |
| 497 | low correlation with all operational dimensions         |     |
| 498 | ( $ r  < 0.13$ ), making it different from other attack |     |
| 499 | types that show a preference for specific opera-        |     |
| 500 | tional patterns. The lack of correlation highlights     |     |
| 501 | its strong implementational flexibility and its inde-   |     |
| 502 | pendence from specific operational constraints.         |     |
| 503 | This characteristic is primarily due to the attack's    |     |
| 504 | low data dependency, relying solely on the inputs       |     |
| 505 | and outputs of the LLMs. Unlike white-box manip-        |     |
| 506 | ulation that prefers perturbed prompts for reliability  |     |
| 507 | ( $r = 0.57$ ), or search-based attacks that typically  |     |
| 508 | require iterative interaction ( $r = 0.61$ ), black-box |     |
| 509 | prefix dependency exploitation adapts well to the       |     |
| 510 | operational constraints discussed in the paper.         |     |
| 511 | Another contributing factor is that this jailbreak      |     |
| 512 | attack category is relatively novel, and a main-        |     |
| 513 | stream operational paradigm has not yet formed.         |     |
| 514 | The diversity within this category can also lead to     |     |
| 515 | a low preference for specific operational patterns.     |     |
| 516 | <b>6 Future Directions</b>                              |     |
| 517 | Through our analysis, we identify several critical      |     |
| 518 | research gaps that remain unexplored or insuffi-        |     |
| 519 | ciently addressed, presenting both challenges and       |     |
| 520 | opportunities for future investigation.                 |     |
| 521 | <b>6.1 Expanding Attack Surface Coverage</b>            |     |
| 522 | <b>Out-of-Distribution Attack Vectors Exploration.</b>  |     |
| 523 | Existing OOD jailbreak attacks mostly rely on           |     |
| 524 | cross-language translation and code transformation,     |     |
| 525 | leaving various other distributional imbalances un-     |     |
| 526 | explored, such as domain-specific jargon or spe-        |     |
| 527 | cialized notation systems. Furthermore, the poten-      |     |
| 528 | tial for multiple OOD attacks to collaborate in a       |     |
| 529 | single attack prompt remains largely unexplored.        |     |
| 530 | Identifying exploitable OOD attack vectors and in-      |     |
| 531 | vestigating their collaborative effects forms a new     |     |
| 532 | research frontier.                                      |     |
| 533 | <b>Service Provider Function Exploitation.</b> Gray-    |     |
| 534 | box jailbreak attacks typically require log-            |     |
| 535 | probability data. As service provider ecosystems        |     |
| 536 | expand, more functions that provide token gener-        |     |
| 537 | ation information may emerge. While fine-tuning         |     |
| 538 | functions have already been exploited, many func-       |     |
| 539 | tional weaknesses remain undiscovered. For exam-        |     |
| 540 | ple, LLM service providers may provide metadata         |     |
|     | such as retrieval confidence scores for Retrieval-      | 541 |
|     | Augmented Generation (RAG) systems or execu-            | 542 |
|     | tion probability rankings for agentic systems. Such     | 543 |
|     | data can become effective sources for gray-box at-      | 544 |
|     | tacks, necessitating rigorous security assessment       | 545 |
|     | of these functions before widespread deployment.        | 546 |
|     | <b>6.2 Systematizing Attack Mechanisms</b>              | 547 |
|     | <b>Implicit Pattern Specification and Optimization.</b> | 548 |
|     | Existing research has identified some exploitable       | 549 |
|     | LLM behavioral patterns, such as role-playing,          | 550 |
|     | persuasive techniques, and authority endorsement.       | 551 |
|     | However, the taxonomy of exploitable behavioral         | 552 |
|     | patterns embedded in LLM training remains in-           | 553 |
|     | complete. Multiple assumed effective manipula-          | 554 |
|     | tion techniques, like emotional manipulation and        | 555 |
|     | social engineering, still need to be specified and de-  | 556 |
|     | veloped into concrete and viable attack paradigms.      | 557 |
|     | Moreover, mapping between specific attack con-          | 558 |
|     | texts and optimal pattern selection lacks systematic    | 559 |
|     | methodology, forming another research frontier.         | 560 |
|     | <b>Clarifying Context Manipulation Mechanisms.</b>      | 561 |
|     | The concrete mechanisms and effective boundaries        | 562 |
|     | of context contamination and context extension          | 563 |
|     | need further clarification. Existing jailbreak meth-    | 564 |
|     | ods often deploy both attacks simultaneously, mak-      | 565 |
|     | ing it difficult to determine the effectiveness of each | 566 |
|     | attack type in isolation. The optimal contamination-    | 567 |
|     | to-extension ratios, the threshold lengths where        | 568 |
|     | extension effects dominate over contamination ef-       | 569 |
|     | fects, and the interaction mechanisms between           | 570 |
|     | them require further investigation to provide deeper    | 571 |
|     | insights into attention dilution attacks.               | 572 |
|     | <b>7 Conclusion</b>                                     | 573 |
|     | This work surveys 43 LLM jailbreak methods from         | 574 |
|     | 36 papers, analyzing them through a vulnerability-      | 575 |
|     | centric framework that examines both exploitation       | 576 |
|     | mechanisms and operational characteristics. Our         | 577 |
|     | correlation analysis reveals systematic patterns be-    | 578 |
|     | tween attack types and operational constraints that     | 579 |
|     | were previously obscured by approach-centric cat-       | 580 |
|     | egorizations. By synthesizing current research, we      | 581 |
|     | identify critical gaps in attack surface coverage,      | 582 |
|     | mechanism systematization, and the effectiveness-       | 583 |
|     | stealthiness trade-off, providing practical guidance    | 584 |
|     | for developing more principled offensive and de-        | 585 |
|     | fensive techniques.                                     | 586 |

|     |  |     |
|-----|--|-----|
| 587 | <b>Ethical Consideration</b>                           |     |
| 588 | This work addresses a critical security challenge in   |     |
| 589 | LLMs while acknowledging the sensitive nature of       |     |
| 590 | jailbreak attack research. We adhere to the follow-    |     |
| 591 | ing ethical principles in conducting this research:    |     |
| 592 | <b>Responsible Disclosure.</b> Our taxonomy is de-     |     |
| 593 | signed to enhance defensive capabilities rather than   |     |
| 594 | facilitate malicious exploits. Through systematic      |     |
| 595 | categorization of vulnerabilities, we seek to em-      |     |
| 596 | power developers and researchers to proactively        |     |
| 597 | strengthen LLM safety. We deliberately focus on        |     |
| 598 | vulnerability taxonomy and attack pattern analysis     |     |
| 599 | while avoiding detailed implementation guides that     |     |
| 600 | could enable jailbreak attacks.                        |     |
| 601 | <b>Broader Impact.</b> This research advances the      |     |
| 602 | safe deployment of LLMs in sensitive domains.          |     |
| 603 | By identifying system vulnerabilities and research     |     |
| 604 | gaps, we provide actionable insights from both         |     |
| 605 | technical and operational perspectives to support      |     |
| 606 | the development of more robust safety mechanisms       |     |
| 607 | that benefit society.                                  |     |
| 608 | <b>Limitations</b>                                     |     |
| 609 | <b>Lack of Standardized Benchmarking</b>               |     |
| 610 | Currently, LLM jailbreaking lacks standardized         |     |
| 611 | testing configurations. Different researchers have     |     |
| 612 | deployed various experimental setups, involving        |     |
| 613 | different target models, resource and access level     |     |
| 614 | requirements, and evaluation criteria for attack suc-  |     |
| 615 | cess. These discrepancies make direct experimental     |     |
| 616 | comparisons difficult and may result in misleading     |     |
| 617 | conclusions. Furthermore, attempting to reproduce      |     |
| 618 | various attack methods under a unified condition       |     |
| 619 | would be methodologically questionable, as each        |     |
| 620 | attack was originally designed and optimized for a     |     |
| 621 | specific context. Given these challenges, we focus     |     |
| 622 | on the construction of a comprehensive taxonomy        |     |
| 623 | rather than conducting large-scale experiments.        |     |
| 624 | <b>Limited Coverage of Defensive Mechanisms</b>        |     |
| 625 | Although our taxonomy identifies vulnerabilities       |     |
| 626 | that defensive mechanisms need to address, we          |     |
| 627 | do not comprehensively survey existing defenses.       |     |
| 628 | This is because existing defensive measures are        |     |
| 629 | highly fragmented, ranging from input filtering and    |     |
| 630 | output monitoring to safety alignment approaches.      |     |
| 631 | Comprehensively surveying these diverse meth-          |     |
| 632 | ods would require a separate taxonomy frame-           |     |
| 633 | work. Furthermore, the opacity of commercial de-       |     |
| 634 | fenses—where many state-of-the-art mechanisms          |     |
|     | remain proprietary—makes systematic categoriza-        | 635 |
|     | tion challenging without access to implementation      | 636 |
|     | details. Instead, we provide a vulnerability-centric   | 637 |
|     | jailbreak attack taxonomy that can serve as a founda-  | 638 |
|     | tion for future defensive taxonomy design.             | 639 |
|     | <b>Limited Quantitative Validation</b>                 | 640 |
|     | Our correlation analysis is based on binary occur-     | 641 |
|     | rence data rather than quantitative metrics. This is   | 642 |
|     | because developing such quantitative criteria faces    | 643 |
|     | great challenges in measurement. There exist no        | 644 |
|     | standards or benchmarks to assess the degree or        | 645 |
|     | proportion of each exploitation mechanism. For         | 646 |
|     | instance, how should we assess what percentage         | 647 |
|     | of an attack’s success derives from implicit pattern   | 648 |
|     | exploitation versus OOD techniques when both are       | 649 |
|     | present? Or how do we quantify the relative con-       | 650 |
|     | tribution of attention dilution versus prefix depen-   | 651 |
|     | dency in a hybrid attack? Therefore, while more        | 652 |
|     | granular measurements might provide additional         | 653 |
|     | insights, the categorical approach we employ of-       | 654 |
|     | fers practical clarity for understanding attack mech-  | 655 |
|     | anisms and their relationships to operational con-     | 656 |
|     | straints without requiring potentially arbitrary quan- | 657 |
|     | titative weightings.                                   | 658 |
|     | <b>Generalizability Across Model Architectures</b>     | 659 |
|     | Our taxonomy is built primarily around                 | 660 |
|     | transformer-based autoregressive LLMs, which           | 661 |
|     | dominate current research and deployment. How-         | 662 |
|     | ever, emerging architectures such as state-space       | 663 |
|     | models (e.g., Mamba) or future paradigms may           | 664 |
|     | introduce novel vulnerabilities not captured           | 665 |
|     | by our framework. Although we believe our              | 666 |
|     | lifecycle-based approach (training-derived vs.         | 667 |
|     | inference-derived) provides some generalizability,     | 668 |
|     | attacks that target architecture-specific features     | 669 |
|     | may require framework extensions.                      | 670 |
|     | <b>Acknowledgement of AI Usage</b>                     | 671 |
|     | In this work, we utilized AI language models for       | 672 |
|     | writing assistance. Specifically, we used Anthropic    | 673 |
|     | Claude Sonnet 4.5 and Google’s Gemini 2.0 to           | 674 |
|     | polish language, improve sentence clarity, and en-     | 675 |
|     | hance readability. These tools were used solely for    | 676 |
|     | linguistic refinement; all conceptual development,     | 677 |
|     | taxonomic frameworks, correlation analyses, litera-    | 678 |
|     | ture review, and scientific insights presented in this | 679 |
|     | work are the original intellectual contributions of    | 680 |
|     | the authors.   | 681 |

682

683  
684  
685  
686687  
688  
689  
690691  
692  
693  
694  
695  
696  
697  
698  
699700  
701  
702703  
704  
705  
706707  
708  
709  
710  
711  
712713  
714  
715  
716717  
718  
719  
720721  
722  
723  
724725  
726  
727  
728  
729  
730731  
732  
733  
734

## References

Francisco Aguilera-Martínez and Fernando Berzal. 2025. [Llm security: Vulnerabilities, attacks, defenses, and countermeasures](#). *Preprint*, arXiv:2505.01177.

Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. 2025. [Jailbreaking leading safety-aligned llms with simple adaptive attacks](#). *Preprint*, arXiv:2404.02151.

Cem Anil, Esin DURMUS, Nina Rimsky, Mrinank Sharma, Joe Benton, Sandipan Kundu, Joshua Batson, Meg Tong, Jesse Mu, Daniel J Ford, Francesco Mosconi, Rajashree Agrawal, Rylan Schaeffer, Naomi Bashkansky, Samuel Svenningsen, Mike Lambert, Ansh Radhakrishnan, Carson Denison, Evan J Hubinger, and 15 others. 2024. [Many-shot jailbreaking](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. [Longformer: The long-document transformer](#). *Preprint*, arXiv:2004.05150.

Emet Bethany, Mazal Bethany, Juan Arturo Nolasco Flores, Sumit Kumar Jha, and Peyman Najafirad. 2024. [Jailbreaking large language models with symbolic mathematics](#). *Preprint*, arXiv:2409.11445.

Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. 2021. [Extracting training data from large language models](#). *Preprint*, arXiv:2012.07805.

Zhiyuan Chang, Mingyang Li, Yi Liu, Junjie Wang, Qing Wang, and Yang Liu. 2024. [Play guessing game with llm: Indirect jailbreak attack with implicit clues](#). *Preprint*, arXiv:2402.09091.

Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2024. [Jailbreaking black box large language models in twenty queries](#). *Preprint*, arXiv:2310.08419.

Xuan Chen, Yuzhou Nie, Wenbo Guo, and Xiangyu Zhang. 2025. [When llm meets drl: Advancing jailbreaking efficiency via drl-guided search](#). *Preprint*, arXiv:2406.08705.

Zhiyu Zoey Chen, Jing Ma, Xinlu Zhang, Nan Hao, An Yan, Armineh Nourbakhsh, Xianjun Yang, Julian McAuley, Linda Petzold, and William Yang Wang. 2024. [A survey on large language models for critical societal domains: Finance, healthcare, and law](#). *Preprint*, arXiv:2405.01769.

Junjie Chu, Yugeng Liu, Ziqing Yang, Xinyue Shen, Michael Backes, and Yang Zhang. 2025a. [Jailbreakradar: Comprehensive assessment of jailbreak attacks against llms](#). *Preprint*, arXiv:2402.05668.

Junjie Chu, Yugeng Liu, Ziqing Yang, Xinyue Shen, Michael Backes, and Yang Zhang. 2025b. [Jailbreakradar: Comprehensive assessment of jailbreak attacks against llms](#). *Preprint*, arXiv:2402.05668. 735  
736  
737  
738

Tianyu Cui, Yanling Wang, Chuanpu Fu, Yong Xiao, Sijia Li, Xinhao Deng, Yunpeng Liu, Qinglin Zhang, Ziyi Qiu, Peiyang Li, Zhixing Tan, Junwu Xiong, Xinyu Kong, Zujie Wen, Ke Xu, and Qi Li. 2024. [Risk taxonomy, mitigation, and assessment benchmarks of large language model systems](#). *Preprint*, arXiv:2401.05778. 739  
740  
741  
742  
743  
744  
745

Ziyan Cui, Ning Li, and Huaikang Zhou. 2025. [Can large language models replace human subjects? a large-scale replication of scenario-based experiments in psychology and management](#). *Preprint*, arXiv:2409.00128. 746  
747  
748  
749  
750

Gelei Deng, Yi Liu, Yuekang Li, Kailong Wang, Ying Zhang, Zefeng Li, Haoyu Wang, Tianwei Zhang, and Yang Liu. 2024a. [Masterkey: Automated jailbreaking of large language model chatbots](#). In *Proceedings 2024 Network and Distributed System Security Symposium*. Internet Society. 751  
752  
753  
754  
755  
756

Yue Deng, Wenxuan Zhang, Sinno Jialin Pan, and Lidong Bing. 2024b. [Multilingual jailbreak challenges in large language models](#). *Preprint*, arXiv:2310.06474. 757  
758  
759  
760

Peng Ding, Jun Kuang, Dan Ma, Xuezhi Cao, Yunsen Xian, Jiajun Chen, and Shujian Huang. 2024. [A wolf in sheep’s clothing: Generalized nested jailbreak prompts can fool large language models easily](#). *Preprint*, arXiv:2311.08268. 761  
762  
763  
764  
765

Matteo Esposito, Francesco Palagiano, Valentina Lenarduzzi, and Davide Taibi. 2025. [On large language models in mission-critical it governance: Are we ready yet?](#) *Preprint*, arXiv:2412.11698. 766  
767  
768  
769

Lang Gao, Jiahui Geng, Xiangliang Zhang, Preslav Nakov, and Xiuying Chen. 2025. [Shaping the safety boundaries: Understanding and defending against jailbreaks in large language models](#). *Preprint*, arXiv:2412.17034. 770  
771  
772  
773  
774

Chuan Guo, Alexandre Sablayrolles, Hervé Jégou, and Douwe Kiela. 2021. [Gradient-based adversarial attacks against text transformers](#). *Preprint*, arXiv:2104.13733. 775  
776  
777  
778

Xingang Guo, Fangxu Yu, Huan Zhang, Lianhui Qin, and Bin Hu. 2024a. [Cold-attack: Jailbreaking llms with stealthiness and controllability](#). *Preprint*, arXiv:2402.08679. 779  
780  
781  
782

Yufei Guo, Muzhe Guo, Juntao Su, Zhou Yang, Mengqiu Zhu, Hongfei Li, Mengyang Qiu, and Shuo Shuo Liu. 2024b. [Bias in large language models: Origin, evaluation, and mitigation](#). *Preprint*, arXiv:2411.10915. 783  
784  
785  
786  
787

|     |  |     |
|-----|--|-----|
| 788 | Zeqing He, Zhibo Wang, Zhixuan Chu, Huiyu Xu, Wenhui Zhang, Qinglong Wang, and Rui Zheng. 2025. Jailbreaklens: Interpreting jailbreak mechanism in the lens of representation and circuit. <i>Preprint</i> , arXiv:2411.11114.   | 843 |
| 789 |  | 844 |
| 790 |  | 845 |
| 791 |  | 846 |
| 792 |  |     |
| 793 | Jie Huang and Kevin Chen-Chuan Chang. 2023. Towards reasoning in large language models: A survey. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1049–1065, Toronto, Canada. Association for Computational Linguistics.   | 847 |
| 794 |  | 848 |
| 795 |  | 849 |
| 796 |  | 850 |
| 797 |  |     |
| 798 | Fengqing Jiang, Zhangchen Xu, Luyao Niu, Zhen Xiang, Bhaskar Ramasubramanian, Bo Li, and Radha Poovendran. 2024. Artprompt: Ascii art-based jailbreak attacks against aligned llms. <i>Preprint</i> , arXiv:2402.11753.  | 851 |
| 799 |  | 852 |
| 800 |  | 853 |
| 801 |  | 854 |
| 802 |  | 855 |
| 803 | Haibo Jin, Ruoxi Chen, Peiyan Zhang, Andy Zhou, Yang Zhang, and Haohan Wang. 2025. Guard: Role-playing to generate natural-language jailbreakings to test guideline adherence of large language models. <i>Preprint</i> , arXiv:2402.03299.  | 856 |
| 804 |  | 857 |
| 805 |  | 858 |
| 806 |  | 859 |
| 807 |  | 860 |
| 808 | Haibo Jin, Leyang Hu, Xinuo Li, Peiyan Zhang, Chonghan Chen, Jun Zhuang, and Haohan Wang. 2024. Jailbreakzoo: Survey, landscapes, and horizons in jailbreaking large language and vision-language models. <i>Preprint</i> , arXiv:2407.01599.  | 861 |
| 809 |  | 862 |
| 810 |  | 863 |
| 811 |  | 864 |
| 812 |  |     |
| 813 | Martin Kuo, Jianyi Zhang, Aolin Ding, Qinsi Wang, Louis DiValentin, Yujia Bao, Wei Wei, Hai Li, and Yiran Chen. 2025. H-cot: Hijacking the chain-of-thought safety reasoning mechanism to jailbreak large reasoning models, including openai o1/o3, deepseek-r1, and gemini 2.0 flash thinking. <i>Preprint</i> , arXiv:2502.12893.                                  | 865 |
| 814 |  | 866 |
| 815 |  | 867 |
| 816 |  | 868 |
| 817 |  | 869 |
| 818 |  | 870 |
| 819 |  | 871 |
| 820 | Andrey Labunets, Nishit V. Pandya, Ashish Hooda, Xiaohan Fu, and Earlene Fernandes. 2025. Fun-tuning: Characterizing the vulnerability of proprietary llms to optimization-based prompt injection attacks via the fine-tuning interface. In <i>2025 IEEE Symposium on Security and Privacy (SP)</i> , page 411–429. IEEE.  | 872 |
| 821 |  | 873 |
| 822 |  | 874 |
| 823 |  | 875 |
| 824 |  | 876 |
| 825 |  |     |
| 826 | Linbao Li, Yannan Liu, Daojing He, and Yu Li. 2025. One model transfer to all: On robust jailbreak prompts generation against llms. <i>Preprint</i> , arXiv:2505.17598.  | 877 |
| 827 |  | 878 |
| 828 |  | 879 |
| 829 |  |     |
| 830 | Xiao Li, Zhuhong Li, Qiongxiu Li, Bingze Lee, Jinghao Cui, and Xiaolin Hu. 2024a. Faster-gcg: Efficient discrete optimization jailbreak attacks against aligned large language models. <i>Preprint</i> , arXiv:2410.15362.   | 880 |
| 831 |  | 881 |
| 832 |  | 882 |
| 833 |  | 883 |
| 834 | Xiaoxia Li, Siyuan Liang, Jiyi Zhang, Han Fang, Aishan Liu, and Ee-Chien Chang. 2024b. Semantic mirror jailbreak: Genetic algorithm based jailbreak prompts against open-source llms. <i>Preprint</i> , arXiv:2402.14872.  | 884 |
| 835 |  | 885 |
| 836 |  | 886 |
| 837 |  | 887 |
| 838 |  | 888 |
| 839 | Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. 2024c. Deepinception: Hypnotize large language model to be jailbreaker. <i>Preprint</i> , arXiv:2311.03191.   | 889 |
| 840 |  | 890 |
| 841 |  | 891 |
| 842 |  | 892 |
|     | Zeyi Liao and Huan Sun. 2024. Amplegcg: Learning a universal and transferable generative model of adversarial suffixes for jailbreaking both open and closed llms. <i>Preprint</i> , arXiv:2404.07921.   | 893 |
|     |  | 894 |
|     |  | 895 |
|     |  | 896 |
|     |  | 897 |
|     | Yuping Lin, Pengfei He, Han Xu, Yue Xing, Makoto Yamada, Hui Liu, and Jiliang Tang. 2024. Towards understanding jailbreak attacks in llms: A representation space analysis. <i>Preprint</i> , arXiv:2406.10794.  |     |
|     |  |     |
|     | Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024a. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.  |     |
|     |  |     |
|     | Tong Liu, Yingjie Zhang, Zhe Zhao, Yinpeng Dong, Guozhu Meng, and Kai Chen. 2024b. Making them ask and answer: Jailbreaking large language models in few queries via disguise and reconstruction. <i>Preprint</i> , arXiv:2402.18104.  |     |
|     |  |     |
|     | Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2024c. Autodan: Generating stealthy jailbreak prompts on aligned large language models. <i>Preprint</i> , arXiv:2310.04451.  |     |
|     |  |     |
|     | Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, and Kailong Wang. 2024d. A hitchhiker’s guide to jailbreaking chatgpt via prompt engineering. <i>Proceedings of the 4th International Workshop on Software Engineering and AI for Data Quality in Cyber-Physical Systems/Internet of Things</i> .                    |     |
|     |  |     |
|     | Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, Kailong Wang, and Yang Liu. 2024e. Jailbreaking chatgpt via prompt engineering: An empirical study. <i>Preprint</i> , arXiv:2305.13860.  |     |
|     |  |     |
|     | Yue Liu, Xiaoxin He, Miao Xiong, Jinlan Fu, Shumin Deng, and Bryan Hooi. 2024f. Flipattack: Jailbreak llms via flipping. <i>Preprint</i> , arXiv:2410.02832.   |     |
|     |  |     |
|     | Huijie Lv, Xiao Wang, Yuansen Zhang, Caishuang Huang, Shihan Dou, Junjie Ye, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024. Codechameleon: Personalized encryption framework for jailbreaking large language models. <i>Preprint</i> , arXiv:2402.16717.   |     |
|     |  |     |
|     | Xingjun Ma, Yifeng Gao, Yixu Wang, Ruofan Wang, Xin Wang, Ye Sun, Yifan Ding, Hengyuan Xu, Yunhao Chen, Yunhan Zhao, Hanxun Huang, Yige Li, Yutao Wu, Jiaming Zhang, Xiang Zheng, Yang Bai, Zuxuan Wu, Xipeng Qiu, Jingfeng Zhang, and 29 others. 2025. Safety at scale: A comprehensive survey of large model and agent safety. <i>Preprint</i> , arXiv:2502.05206. |     |
|     |  |     |
|     | Wenlong Meng, Fan Zhang, Wendao Yao, Zhenyuan Guo, Yuwei Li, Chengkun Wei, and Wenzhi Chen. 2025. Dialogue injection attack: Jailbreaking llms through context manipulation. <i>Preprint</i> , arXiv:2503.08195.   |     |

|     |  |      |
|-----|--|------|
| 898 | Anselm Paulus, Arman Zharmagambetov, Chuan Guo, Brandon Amos, and Yuandong Tian. 2025. <a href="#">Advprompter: Fast adaptive adversarial prompting for llms</a> . <i>Preprint</i> , arXiv:2404.16873.   | 955  |
| 899 |  | 956  |
| 900 |  | 957  |
| 901 |  | 958  |
| 902 | Matthew Renze and Erhan Guven. 2024. <a href="#">The benefits of a concise chain of thought on problem-solving in large language models</a> . In <i>2024 2nd International Conference on Foundation and Large Language Models (FLLM)</i> , page 476–483. IEEE.   | 959  |
| 903 |  | 960  |
| 904 |  | 961  |
| 905 |  |      |
| 906 |  |      |
| 907 | Sippo Rossi, Alisia Marianne Michel, Raghava Rao Mukkamala, and Jason Bennett Thatcher. 2024. <a href="#">An early categorization of prompt injection attacks on large language models</a> . <i>Preprint</i> , arXiv:2402.00898.   | 962  |
| 908 |  | 963  |
| 909 |  | 964  |
| 910 |  | 965  |
| 911 | Mikayel Samvelyan, Sharath Chandra Raparthy, Andrei Lupu, Eric Hambro, Aram H. Markosyan, Manish Bhatt, Yuning Mao, Minqi Jiang, Jack Parker-Holder, Jakob Foerster, Tim Rocktäschel, and Roberta Raileanu. 2024. <a href="#">Rainbow teaming: Open-ended generation of diverse adversarial prompts</a> . <i>Preprint</i> , arXiv:2402.16822.                                      | 966  |
| 912 |  | 967  |
| 913 |  | 968  |
| 914 |  | 969  |
| 915 |  | 970  |
| 916 |  | 971  |
| 917 |  | 972  |
| 918 | Sander Schulhoff, Jeremy Pinto, Anaam Khan, Louis-François Bouchard, Chenglei Si, Svetlana Anati, Valen Tagliabue, Anson Liu Kost, Christopher Carnahan, and Jordan Boyd-Graber. 2024. <a href="#">Ignore this title and hackaprompt: Exposing systemic vulnerabilities of llms through a global scale prompt hacking competition</a> . <i>Preprint</i> , arXiv:2311.16119.        | 973  |
| 919 |  | 974  |
| 920 |  | 975  |
| 921 |  | 976  |
| 922 |  |      |
| 923 |  |      |
| 924 |  |      |
| 925 | Rusheb Shah, Quentin Feuillade-Montixi, Soroush Pour, Arush Tagade, Stephen Casper, and Javier Rando. 2023. <a href="#">Scalable and transferable black-box jailbreaks for language models via persona modulation</a> . <i>Preprint</i> , arXiv:2311.03348.  | 977  |
| 926 |  | 978  |
| 927 |  | 979  |
| 928 |  | 980  |
| 929 |  |      |
| 930 | Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. 2023. <a href="#">Large language models can be easily distracted by irrelevant context</a> . <i>Preprint</i> , arXiv:2302.00093.  | 981  |
| 931 |  | 982  |
| 932 |  | 983  |
| 933 |  | 984  |
| 934 |  |      |
| 935 | Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. <a href="#">AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4222–4235, Online. Association for Computational Linguistics. | 985  |
| 936 |  | 986  |
| 937 |  | 987  |
| 938 |  | 988  |
| 939 |  | 989  |
| 940 |  |      |
| 941 |  |      |
| 942 | Yu Tian, Xiao Yang, Jingyuan Zhang, Yinpeng Dong, and Hang Su. 2024. <a href="#">Evil geniuses: Delving into the safety of llm-based agents</a> . <i>Preprint</i> , arXiv:2311.11855.  | 990  |
| 943 |  | 991  |
| 944 |  | 992  |
| 945 |  | 993  |
| 946 | Zeming Wei, Yifei Wang, Ang Li, Yichuan Mo, and Yisen Wang. 2024. <a href="#">Jailbreak and guard aligned language models with only few in-context demonstrations</a> . <i>Preprint</i> , arXiv:2310.06387.  | 994  |
| 947 |  | 995  |
| 948 |  | 996  |
| 949 |  | 997  |
| 950 | Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. 2023. <a href="#">Hard prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery</a> . <i>Preprint</i> , arXiv:2302.03668.  | 998  |
| 951 |  | 999  |
| 952 |  | 1000 |
| 953 |  | 1001 |
| 954 |  |      |
|     | Tianyi Wu, Zhiwei Xue, Yue Liu, Jiaheng Zhang, Bryan Hooi, and See-Kiong Ng. 2025. <a href="#">Geneshift: Impact of different scenario shift on jailbreaking llm</a> . <i>Preprint</i> , arXiv:2504.08104.   | 1002 |
|     |  | 1003 |
|     |  | 1004 |
|     |  | 1005 |
|     | Zeguan Xiao, Yan Yang, Guanhua Chen, and Yun Chen. 2024. <a href="#">Distract large language models for automatic jailbreak attack</a> . <i>Preprint</i> , arXiv:2403.08424.   |      |
|     |  |      |
|     | Zihao Xu, Yi Liu, Gelei Deng, Yuekang Li, and Stjepan Picek. 2024. <a href="#">A comprehensive study of jailbreak attack versus defense for large language models</a> . <i>Preprint</i> , arXiv:2402.13457.  |      |
|     |  |      |
|     | Lixiang Yan, Lele Sha, Linxuan Zhao, Yuheng Li, Roberto Martinez-Maldonado, Guanliang Chen, Xinyu Li, Yueqiao Jin, and Dragan Gašević. 2023. <a href="#">Practical and ethical challenges of large language models in education: A systematic scoping review</a> . <i>British Journal of Educational Technology</i> , 55(1):90–112.  |      |
|     |  |      |
|     | Sibo Yi, Yule Liu, Zhen Sun, Tianshuo Cong, Xinlei He, Jiaying Song, Ke Xu, and Qi Li. 2024. <a href="#">Jailbreak attacks and defenses against large language models: A survey</a> . <i>Preprint</i> , arXiv:2407.04295.  |      |
|     |  |      |
|     | Jiahao Yu, Xingwei Lin, Zheng Yu, and Xinyu Xing. 2024. <a href="#">Gptfuzzer: Red teaming large language models with auto-generated jailbreak prompts</a> . <i>Preprint</i> , arXiv:2309.10253.   |      |
|     |  |      |
|     | Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. 2024. <a href="#">Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher</a> . <i>Preprint</i> , arXiv:2308.06463.  |      |
|     |  |      |
|     | Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. 2024. <a href="#">How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms</a> . <i>Preprint</i> , arXiv:2401.06373.  |      |
|     |  |      |
|     | Xiaosen Zheng, Tianyu Pang, Chao Du, Qian Liu, Jing Jiang, and Min Lin. 2024. <a href="#">Improved few-shot jailbreaking can circumvent aligned language models and their defenses</a> . <i>Preprint</i> , arXiv:2406.01288.   |      |
|     |  |      |
|     | Yuyang Zhou, Guang Cheng, Kang Du, Zihan Chen, and Yuyu Zhao. 2025. <a href="#">Toward intelligent and secure cloud: Large language model empowered proactive defense</a> . <i>Preprint</i> , arXiv:2412.21051.  |      |
|     |  |      |
|     | Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023. <a href="#">Universal and transferable adversarial attacks on aligned language models</a> . <i>Preprint</i> , arXiv:2307.15043.   |      |
|     |  |      |
|     | <b>A Jailbreak Formalization</b>   |      |
|     | To formalize the process of LLM jailbreak, we first have to identify the elements involved. Let $x$ denote the initial harmful prompt that would typically   |      |

be rejected by the LLM, and let  $f(x)$  be a transformation function that converts the original prompt into a modified version  $x'$ .

$$x' = f(x) \quad (1)$$

Let LLM represent the target language model, and  $y$  be the generated response to the manipulated prompt.

$$y = \text{model}(x') \quad (2)$$

Let JUDGE be an evaluation function that returns True if the response  $y$  violates the safety objective  $G$  (indicating harmful content generation), and False otherwise.

$$R = \text{JUDGE}(y, G) | R \in \{True, False\} \quad (3)$$

This means that a jailbreak attack is considered successful if the protective mechanisms of the target LLM are circumvented (i.e., the target LLM and its defensive mechanisms fail to identify the malicious intention of the prompt), and harmful content is generated (i.e., the result of JUDGE is True).

## B Jailbreak Analysis according to Technical Dimensions

### B.1 Training-derived attacks

Jailbreak attacks that fall into this category all target the vulnerabilities embedded in the training phase, i.e., the pre-training and safety alignment phase. These attacks can be assigned to the following types:

#### B.1.1 Out-of-Distribution(OOD) Attack

Distributional imbalances across multiple dimensions in the training datasets for LLMs have been observed (Guo et al., 2024b). These distributional imbalances result in insufficiency in the model’s safety alignment in various domains and contexts. As a result, models exhibit reduced safety awareness in these scenarios, and this vulnerability resulted in the success of Out-of-Distribution (OOD) attacks.

We define the out-of-distribution (OOD) attack function as:

$$f_{OOD}(x) = O^*(x) \quad (4)$$

$$\text{s.t. } O^* = \arg \max_{O \in \mathcal{O}} \mathcal{L}(\text{model}(O(x)), y_{\text{harmful}}) \cdot (1 - P_Q(O(x), \delta))$$

- $x$  represents the original harmful request
- $O \in \mathcal{O}$  represents transformation functions from the OOD manipulation space that modify input to maximize distributional divergence
- $y_{\text{harmful}}$  denotes the target harmful output
- $\mathcal{L}(\text{model}(O(x)), y_{\text{harmful}})$  measures the model’s tendency to generate harmful content when given transformed input  $O(x)$
- $P_Q(O(x), \delta)$  denotes the probability density function of the training distribution  $Q$  evaluated at transformed input  $O(x)$ ;  $\delta$  represents perturbation parameters that exploit identified distributional weaknesses/gaps in training data  $Q$ .

The key is to find the transformation  $O^*$  that maximizes the correspondence of harmful output and meanwhile minimizing the likelihood of content appearing frequency in training data (exploiting OOD regions where safety training is sparse) according to  $\delta$ .

Although the training data for most commercial LLMs is inaccessible, attackers can still get relative information by analyzing output patterns and making theoretical assumptions about the training data (Carlini et al., 2021). Through this analysis, they can design attack prompts that target these domains with insufficient data. OOD attacks can be divided into two distinct types, depending on the  $\delta$  they deployed:

**Translation-based:** Researchers have found that most existing LLMs safety works focus on English contexts. Challenges still exist in safety alignment in multilingual (e.g. the mixture of Chinese and English) and low-resource language environments (e.g. Arabic, German, etc). This imbalance can be formalized as  $\delta \in \mathcal{S}_{\text{linguistic}}$  (language space).

ReNeLLM (Ding et al., 2024) represents a group of OOD attacks that leverage the multilingual environments. It employs a two-stage translation approach to partially convert harmful prompts from English to Chinese as part of its attack pipeline, in order to create an exploitable multilingual environment.

MultiJail (Deng et al., 2024b) represents the other group of OOD attack methods that exploit the low-resource language environment. With native speakers translating harmful prompts from English to 9 distinct languages, the researchers have constructed a jailbreak dataset. They also developed an automated translation module that can traverse all 9 languages to execute attacks.

**Non-Translation-based:** Besides languages, researchers have also found distributional flaws in other domains. This can be formulated as  $\delta \in \mathcal{S}_{\text{representational}}$  (encoding/structural space).

ArtPrompt (Jiang et al., 2024) transforms the harmful content in the prompt into ASCII art form without changing its position in the sentences. The success of this attack exposed the weakness of LLM safety mechanisms in detecting harmful content in ASCII art form, since the target LLMs themselves can still understand the malicious purpose of the prompt.

DRA (Liu et al., 2024b) identifies a theoretical distributional gap in LLM training data: LLMs are exposed to harmful instructions at a higher frequency in query contexts than in completion contexts in safety alignment phases. This gap results in reduced safety awareness when processing ostensibly self-generated harmful content.

These attacks that target the distributional flaws in training data reveal a critical challenge: increasing data coverage for pre-training and the safety alignment phase requires more computational resources and extended training time, and could suffer from diminishing marginal utility. Therefore, the development of more adaptive and generalizable safety mechanisms that can remain robust against harmful prompts regardless of data dependencies has become a must.

### B.1.2 Implicit Pattern Attack

OOD attacks originate from domains that are low in data frequency in training data. Conversely, Implicit Pattern attacks originate from the domains with high data frequency. It is this high data frequency that creates shared vulnerabilities among most existing LLMs.

These shared vulnerabilities can be traced back to implicit patterns embedded during training phases through similarities in training methodologies and data sources across different LLMs. According to research by Cui et al. (2025), responses from LLMs show strong similarities to human behavioral patterns. It is these behavioral simi-

larities that create exploitable vulnerabilities that are shared by multiple LLMs, as demonstrated by Liu et al. (2024e), who have documented several widely applicable attack patterns, including persona adoption and psychological manipulation, etc. These strategies have already been the foundations of various jailbreak methods. However, the mapping between specific attack contexts and their optimal pattern selection is still largely implicit. Meanwhile, current behavioral patterns are identified through analyzing the LLM output; so far, there exists no systematic method for human-like vulnerability categorization, leaving numerous patterns undiscovered.

We define Implicit Pattern attacks as exploiting shared vulnerabilities across LLMs through behavioral mimicry patterns. Let  $\Psi = \{\psi_1, \psi_2, \dots, \psi_k\}$  represent the set of behavioral patterns embedded during training, where each pattern  $\psi_i$  corresponds to exploitable similarities in human behaviors. The implicit pattern attack function can be formulated as:

$$f_{\text{implicit}}(x) = g^*(x) \quad (5)$$

$$\text{where } g^* = \arg \max_{g \in \mathcal{G}} \sum_{i=1}^k \quad (6)$$

$$\mathcal{L}(\text{model}(g(x)), y_{\text{harmful}}) \cdot w_i(x) \cdot S(g(x), \psi_i) \quad (7)$$

- $x$  represents the original harmful request
- $\Psi = \{\psi_1, \psi_2, \dots, \psi_k\}$  represents the set of behavioral patterns embedded during training;  $\psi_i$  corresponds to exploitable similarities in human behaviors for pattern  $i$
- $g \in \mathcal{G}$  represents transformation functions from the attack framework space (persona adoption, psychological manipulation, other undiscovered patterns, etc.)
- $y_{\text{harmful}}$  denotes the target harmful output
- $w_i(x)$  represents the compatibility weight of pattern  $\psi_i$  with respect to the original prompt  $x$ , indicating how appropriate the pattern is for the specific type of harmful request  $\psi_i$ . The weights indicate the importance of corresponding behavioral patterns in the attack strategy
- $S(g(x), \psi_i)$  quantifies the similarity between the transformed attack prompt  $g(x)$  and human behavioral pattern  $\psi_i$

1190 •  $\mathcal{L}(\text{model}(g(x)), y_{\text{harmful}})$  ensures the model's  
1191 response corresponds to the original harmful  
1192 intent from input  $x$

1193 •  $k$  represents the total number of implicit behav-  
1194 ioral patterns in the set  $\Psi$

1195 Optimizing  $g^*$  means maximizing weighted be-  
1196 havioral similarity across all patterns and maintain-  
1197 ing corresponding harmful output.

1198 These implicit pattern attacks can be categorized  
1199 into two approaches based on how they acquire a  
1200 suitable  $S(g(x), \Psi)$  matrix:

1201 **Search-Based:** Search-based approaches get the  
1202 optimal  $S(g(x), \Psi)$  matrix by searching through  
1203 the  $\Psi$  space, i.e., searching for the combination of  
1204 behavioral patterns that can achieve jailbreaks.

1205 AutoDAN (Liu et al., 2024c), Geneshift (Wu  
1206 et al., 2025), and SMJ (Li et al., 2024b) represent  
1207 methods that employ genetic algorithms for the  
1208 searching process.

1209 GPTFUZZER (Yu et al., 2024) and PAIR (Chao  
1210 et al., 2024) represent methods that apply broad  
1211 search strategies for every single query.

1212 Rainbowteaming (Samvelyan et al., 2024) redef-  
1213 ines the jailbreak challenge as a Quality-Diversity  
1214 (QD) search optimization problem and expand the  
1215 search process into multi-dimensional spaces.

1216 All these methods can generate sophisticated  
1217 jailbreak prompts embedded with multiple behav-  
1218 ior patterns according to different contexts while  
1219 requiring moderate computational resources.

1220 **Learning-Based:** Learning-based approaches  
1221 acquire the optimal  $S(g(x), \Psi)$  matrix by learning  
1222 from existing contexts-to- $S(g(x), \Psi)$  mapping of  
1223 existing successful jailbreak examples.

1224 ICA (Wei et al., 2024), MSJ (Anil et al., 2024),  
1225 and I-SFJ (Zheng et al., 2024) utilize the in-context  
1226 learning ability of LLMs to simulate the behavioral  
1227 patterns from the examples presented directly in  
1228 the prompts.

1229 RLbreaker (Chen et al., 2025) deploys a rein-  
1230 forcement learning strategy to train the optimizer  
1231 in the prompt generation pipeline.

1232 Arrattack (Li et al., 2025) and Adv-  
1233 prompter (Paulus et al., 2025) fine-tune  
1234 LLMs to learn behavioral patterns from successful  
1235 jailbreak datasets.

1236 These methods can achieve consistent and trans-  
1237 ferable jailbreak performance, and could even learn  
1238 how to circumvent recent defensive mechanisms if  
1239 the dataset is up to date.

The success of Implicit Pattern attacks exposes  
a critical safety challenge by exploiting the be-  
havioral patterns embedded within LLMs during  
training. The effectiveness of the search-based and  
learning-based approaches shows that both the com-  
prehensive exploration of attack spaces and pattern  
replication can be efficient in generating effective  
jailbreak prompts.

## B.2 Inference-derived attacks

Jailbreak attacks that fall into this category all  
target the vulnerabilities existing in the inference  
phase, where the models process the received data  
and generate outputs. These attacks can be as-  
signed to the following categories:

### B.2.1 Attention Dilution Attack

Attention dilution attacks refer to attacks that ex-  
ploit the weaknesses in the attention mechanisms  
of transformer-based models. These attacks modify  
the input prompts in a way that can misdirect or  
overload the model's attention to reduce the pos-  
sibility of the malicious intention being detected.  
The key reason for the success of these attacks  
is that transformer models only possess limited  
attention resources to allocate across tokens, and  
this allocation can be manipulated through inten-  
tional prompt modification, i.e., attention would  
be diluted in this process. This dilution prevents  
the models' safety mechanisms from focusing on  
potentially harmful contents and forces them to  
eventually generate harmful contents.

Liu et al. (2024a) have shown that even LLMs  
explicitly designed for long-context processing ex-  
perience a notable performance reduction within  
the middle segments of extended inputs. The ori-  
gin of this vulnerability is the attention mechanism  
itself. The Multi-Head Attention (MHA) mecha-  
nism can build a comprehensive relationship ma-  
trix between every token in the input, but the ex-  
plosive, quadratic increase in required computa-  
tional resources based on input length limits the  
models' ability to process long contexts. Although  
researchers have developed sparse attention mecha-  
nisms for long-context processing, they are capable  
of considering relationships between only a limited  
number of tokens, and thus, they have only limited  
comprehensive context cognition abilities. Attack-  
ers can exploit these potential weaknesses of target  
models through prompt manipulation.

We define Attention Dilution attacks as exploit-  
ing the limited attention allocation in transformer-

based models to bypass safety mechanisms by strategically manipulating input prompts to turn the model’s attention away from harmful content. Let  $A(x, i)$  represent the attention weight allocated to token  $i$  in input sequence  $x$ , subject to the constraint  $\sum_{i=1}^{|x|} A(x, i) = 1$ . The attention dilution attack function can be formulated as:

$$f_{\text{dilution}}(x) = T^*(x)$$

$$\text{s.t. } T^* = \arg \max_{T \in \mathcal{T}} \mathcal{L}(\text{model}(T(x)), y_{\text{harmful}}) \cdot D(T(x))$$

- $x$  represents the original harmful prompt
- $T \in \mathcal{T}$  represents transformation functions from the dilution strategy space that modify the input to dilute attention
- $y_{\text{harmful}}$  denotes the target harmful output
- $\mathcal{L}(\text{model}(T(x)), y_{\text{harmful}})$  measures attack success likelihood for the transformed input
- $D(T(x))$  represents the dilution factor that quantifies how effectively attention is diverted from harmful content

$D(T(x))$  is defined as:

$$D(T(x)) = \frac{|T(x)|}{|x|} \cdot \frac{\sum_{i \in I_{\text{distract}}} A(T(x), i)}{\sum_{i \in I_{\text{harm}}} A(T(x), i)} \quad (8)$$

where  $I_{\text{distract}}$  represents indices of distracting tokens (padding, benign context, encoding artifacts);  $I_{\text{harm}}$  denotes the indices of harmful content tokens in the transformed input  $T(x)$ .  $A(x, i)$  represents the attention weight allocated to token  $i$  in input sequence  $x$ , subject to constraint  $\sum_{i=1}^{|x|} A(x, i) = 1$

The attacks succeed when attention is maximally diverted from harmful content while maintaining the ability to generate corresponding harmful output.

Attacks of this type can be divided into the following two categories, based on which part of  $D(T(x))$  dominates:

**Context-Extension-Based:** The phenomenon of important information capturing ability decay in long sequences, sometimes referred to as the Long-tail Attention Decay effect, has already

been discussed in multiple works, such as Longformer (Beltagy et al., 2020). This work shows that Transformer architectures are inherently unable to process long sequences: Even optimized implementations could only retain self-attention operations that scale linearly with sequence length. When sequences become sufficiently long as the harmful content remains short, the term  $\frac{|T(x)|}{|x|}$  dominates in  $D(T(x))$ . This results in an imbalance in models’ attention allocation between harmful content and other content.

Such an imbalance creates opportunities for context-extension-based attacks. MSJ (Anil et al., 2024) enlarges the scale of demonstration sets and creates an extremely long input context; Mathprompt (Bethany et al., 2024) also extends its input to a much longer context, but with an alternative approach: It embeds/transforms the original harmful prompt into a long and complex task. Both approaches have successfully reduced the percentage of harmful content in the full context.

These attacks leverage the attention allocation constraints of Transformer architectures to circumvent safety mechanisms of the target LLMs via context extension.

**Context-Contamination-Based:** Research (Shi et al., 2023) has showed that LLMs’ performance degrades when input prompts contain information that is irrelevant to their primary intentions. This observation reveals a new attack vector: context contamination. Attackers can intentionally embed the harmful content into irrelevant and seemingly harmless contexts, where  $\frac{\sum_{i \in I_{\text{distract}}} A(T(x), i)}{\sum_{i \in I_{\text{harm}}} A(T(x), i)}$  dominates in  $D(T(x))$ . According to our observations, contaminations appear frequently in role-playing scenarios, task completion scenarios, and special token injection. And meanwhile, they are among the most exploitable scenarios.

Persona modulation (Shah et al., 2023) is a typical method that implements role-playing in the attack pipeline. It creates personas that can potentially generate harmful outputs and formats them as a prompt generation framework.

Regarding task completion scenarios, cipher and decipher operations are widely deployed. Self-Cipher, CipherChat (Yuan et al., 2024), and FLI-Pattack (Liu et al., 2024f) employ encoding techniques to mask harmful content and embed them within seemingly harmless decoding tasks to circumvent LLM safety mechanisms while retain-

ing the malicious intention of the prompt. Other task completion scenarios, including mathematical problem-solving scenarios exploited in Mathprompt (Bethany et al., 2024) and template completion tasks exploited in ReNeLLM (Ding et al., 2024), similarly exploit this attentional vulnerability by distracting the LLM with specific tasks.

Another branch is special token injection. DIA I (Meng et al., 2025) and I-SFJ (Zheng et al., 2024) embed harmful content within deceptive format, especially using the special tokens that could strongly distract the model (e.g., `</s>`, `<assistant>`, or other system tokens). These attacks create an attentional camouflage with injected special tokens. These injected tokens mislead the model’s attention and cause detection mechanisms to fail.

These context contamination techniques demonstrate how safety alignment mechanisms can be circumvented through attention distraction. The harmful content could be hidden behind the primary tasks; this could lead to low attention allocation to harmful content and thus create exploitable vulnerabilities.

## B.2.2 Prefix Dependency Exploitation

Previously introduced Attention Dilution attacks focus on manipulating the attention mechanism during the input processing phase. In contrast, Prefix Dependency Exploitation targets the phase of output generation. Due to the inherent conditional probability dependency characteristic of the autoregressive architectures, LLMs themselves are usually unable to look back and correct their generated output while the generation process is still running. This is further confirmed by the work of Renze et al. (2024), which shows that LLMs require explicit instruction guidance to perform self-reflection processes, indicating that traditional LLMs are unable to autonomously self-reflect even after general pre-training and safety alignment. The models cannot independently “step back” to examine whether their behavior is appropriate during the generation process itself, lacking global safety assessment capabilities. As a consequence, attackers could easily manipulate the LLM outputs by guiding the generation of a benign beginning at first and then harmful content later. Although the appearance of reasoning models possesses reflecting ability, this structural vulnerability can still be exploited.

We define Prefix Dependency Exploitation attacks as exploiting the autoregressive generation mechanism of transformer-based models to ma-

nipulate the conditional probability dependencies during the generation phase, forcing models toward problematic outputs through seemingly benign introductory sequences. Let  $P(x_{t+1}|x_1, x_2, \dots, x_t)$  represent the conditional probability of generating token  $x_{t+1}$  given the sequence prefix  $x_1, x_2, \dots, x_t$ .

The prefix dependency manipulation function  $S(x)$  takes input  $x$  and modifies it based on the first few token probabilities of  $\text{model}(x)$ :

$$S(x) = x \oplus \Delta(x) \quad (9)$$

where  $\Delta(x)$  is computed based on the initial token probability distribution:

$$\begin{aligned} \Delta(x) = & f_{modify}(P(x_1|\text{model}(x)), \\ & P(x_2|x_1, \text{model}(x)), \\ & \dots, P(x_k|x_1 \dots x_{k-1}, \text{model}(x))) \end{aligned} \quad (10)$$

for the first  $k$  tokens in the model’s response. The prefix dependency exploitation attack function can then be formulated as:

$$f_{prefix}(x) = S^*(x) \quad (11)$$

$$\begin{aligned} \text{s.t. } S^* = & \arg \max_{S \in \mathcal{S}} \mathcal{L}(\text{model}(S(x)), y_{harmful}) \\ & \cdot C(S(x)) \cdot (1 - R(S(x))) \end{aligned} \quad (12)$$

where:

- $x$  represents the original harmful request
- $S \in \mathcal{S}$  represents prefix manipulation functions that modify  $x$  based on initial output token probabilities
- $y_{harmful}$  denotes the target harmful output
- $\mathcal{L}(\text{model}(S(x)), y_{harmful})$  measures attack success likelihood
- $C(S(x))$  represents the cumulative dependency factor
- $R(S(x))$  represents the retrospective correction capability (ideally minimized)

The cumulative dependency factor is defined as:

$$C(S(x)) = \prod_{t=1}^{|S(x)|} \frac{P(x_{t+1}^{compliant} | x_1^{benign}, \dots, x_t^{benign})}{P(x_{t+1}^{refused} | x_1^{benign}, \dots, x_t^{benign})} \quad (13)$$

where  $x_t^{benign}$  represents tokens in the manipulated benign prefix and  $x_{t+1}^{compliant}$  represents tokens leading toward compliance with the harmful request.

The self-reflection capability is modeled as:

$$R(S(x)) = \sum_{t=k}^{|S(x)|} P(\text{self reflection} | x_1, \dots, x_t) \cdot I(t > t_{critical}) \quad (14)$$

where  $k$  is the position where harmful intent becomes apparent,  $t_{critical}$  is the point beyond which the self-reflection mechanism becomes weak, and  $I(\cdot)$  is an indicator function.

The jailbreak attacks can be divided into three categories, based on how they acquire the data essential for prefix manipulation:

#### White-Box Manipulation:

Jailbreak attacks belonging to this category acquire essential data through direct access to the model’s parameters.

Greedy Coordinate Gradient (GCG) (Zou et al., 2023) represents the foundational method of this category. It searches for the optimal suffix that maximizes the probability of the first generated token to be an affirmative response like “Sure”, with the assistance of gradient information. Gradients can provide information directly related to token generation probabilities, and can thus transparently guide the suffix optimization process.

Based on GCG, AmpleGCG (Liao and Sun, 2024) expands the search space for suffixes. Instead of using only the best suffix candidate, it considers multiple candidates with high performance to avoid being trapped in local optima, thereby improving attack robustness and ASR. Faster-GCG (Li et al., 2024a) goes a step further. It adds extra terms into the loss function to apply more constraints on the suffix searching process to better customize the attack pipeline according to the scenarios.

In summary, white box manipulation represents a group of attacks that conduct transparent optimization to modify the input prompts that can reliably manipulate language model outputs.

#### Gray-Box Manipulation:

Although White-Box Manipulation provides great explainability and transparency of the attack mechanism, its deployment requires full model access. This restricts its application in most practical situations, where target models do not support full model access. However, the researchers have observed that even closed-source commercial LLMs may expose partial internal information that is still effective in guiding the suffix optimization process. This is a typical application scenario of Gray-Box Manipulation, where malicious users violate the limited data allowed by the LLM vendors.

Some gray-box data are directly provided by the service provider. Adaptive attack (Andriushchenko et al., 2025) replaces the gradient information with log probability data, which is provided by certain closed-source models like GPT-4. These data can work similarly to gradient information in the suffix optimization process, proving the value of gray-box data in the generation manipulation process.

Another approach to retrieve gray-box data is to exploit the interfaces provided by the LLM service providers. A representative example is Fun-tuning (Labunets et al., 2025), where researchers derive data from the fine-tuning API to simulate log probability information and effectively circumvent the need for direct model access.

Gray-Box Manipulation bridges the gap between white-box and black-box approaches by exploiting partially accessible information from closed-source LLMs. Typical Gray-Box Manipulations usually involve exploiting intermediate data and platform-specific interfaces, and they reveal the key fact that any information related to the token generation process can be a potential target for exploitation.

#### Black-Box Manipulation:

White-box and gray-box attacks require extra data beyond the model’s outputs, yet in most strict situations, only the model responses are available to users. Meanwhile, as jailbreak defensive mechanisms develop, exploitable data types are gradually decreasing. In contrast, model inputs and outputs required by Black-Box Manipulation remain and will always remain accessible over time. Furthermore, considering that nonsensical suffixes are often the optimal result of white-box or gray-box manipulation, the stealthiness of these attacks is relatively lower compared to Black-Box Manipulation. Consequently, Black-Box Manipulation has gained importance recently.

Dialogue Injection Attack (DIA) (Meng et al.,

2025) represents one of the most advanced black-box manipulation methods. The DIA I approach (Meng et al., 2025) directly orders the LLM to start from a predefined, affirmative beginning; DIA II applies a two-phase strategy: an initial benign query followed by a malicious query.

H-CoT (Kuo et al., 2025) represents the adaptation of the black-box paradigm to the research context of reasoning models. It not only focuses on generating benign prefixes but also attempts to circumvent the subsequent justification phase within the reasoning model’s context by manipulating the generation of specific tokens in the middle of the output sequences.

Black-box Manipulation relies solely on prompt engineering and input-output analysis. This technique will remain technically viable in the future as it only requires access to the model’s input and output. These methods create deceptive prompts without any nonsensical suffixes. This is an appendix.

## C Jailbreak Analysis according to Operational Dimensions

### C.1 Interaction Methods

There are mainly three ways of interaction between LLM and the attackers. Each method has its own strengths and shortcomings:

**One-Shot:** This approach aims to successfully jailbreak target LLMs using a single attack query. This means that every single malicious prompt has only one chance to conduct the attack, and no feedback information is available for its optimization. This limitation requires researchers to develop generation strategies that could generate high-performance prompts efficiently, and these strategies fall into the following three main categories.

The first category is the surrogate model. Surrogate model methods use substitute LLMs as proxies for target models during attack development (Liu et al., 2024c; Wu et al., 2025; Li et al., 2024b). Attackers assume that these surrogate models behave like the target LLM, and use them to evaluate prompt effectiveness and give feedback when refusals happen. During the prompt generation, the surrogate model provides fitness scores that can guide the optimization process through techniques like genetic algorithms. However, since surrogate models cannot perfectly replicate target model behaviors, the problem of transfer learning appears

– prompts that are optimized for surrogate models may not work equally well on target models.

The second category is the universal template. Universal template approaches create universally applicable prompt templates that are designed for exploiting common weaknesses observed in multiple LLMs (Li et al., 2024c; Bethany et al., 2024). These templates usually combine various attack strategies, such as role-playing scenarios, emotional manipulation, nested instructions, fictional scenarios, or encoding methods to disguise the harmful prompts and form an attack prompt template that contains various attack vectors. Such templates can cover the weakness of multiple LLMs and thus generalize. Developing such templates requires the researchers to have a broad understanding of how different LLMs behave and what common weaknesses are embedded in their safety mechanisms.

The third category is fine-tuned models. Attackers can fine-tune models to automatically generate jailbreak prompts with existing successful jailbreak examples (Liao and Sun, 2024; Li et al., 2025; Paulus et al., 2025). These approaches skip the prompt optimization phase by learning direct mappings from attack objectives to effective prompts, and can produce effective attack prompts in a single forward pass, improving the generation efficiency.

The one-shot paradigm has several notable advantages. Most importantly, it provides stealth and success rate, since the one-shot configuration can circumvent the detection mechanisms that monitor for repeated attack attempts. Moreover, one-shot attacks are also effective in practical scenarios where attackers would usually face the problem of access limitation or rate limitation set by the service provider. Additionally, successful one-shot methods show high generalizability, since they are designed to be capable of working across diverse models without extra modification. Despite so many advantages, this approach still faces some limitations. Its nature of lacking iterative feedback prevents the attack prompt from becoming optimal for different situations. And this may result in lower ASR compared to iterative attacks that can optimize the attack prompt based on model responses.

**Iterative:** Iterative attacks attempt to jailbreak target LLMs by conducting multi-turn queries and performing in-loop optimization of the prompt according to the feedback from the target LLMs. According to the type of feedback retrieved from the

target model, these attacks could be divided into two categories:

First category directly uses the target model’s response to the attack prompt as feedback to guide prompt optimization (Ding et al., 2024; Shah et al., 2023; Jin et al., 2025). These black-box methods derive optimizing strategies from the information given by the target model in their refusal responses and adjust the prompt accordingly, and the most widely used optimization approaches are heuristic search and RL. This allows for the exploration of diverse attack possibilities without requiring internal model access. However, this searching space is constrained by the observable outputs, and in extreme situations, like generation interruption, no reasonable optimization could be made. Besides, analyzing the model’s response is a typical natural language process (NLP) task, which is more difficult than analyzing data like gradients and therefore usually requires extra computational effort.

The second category acquires internal model parameters either in or after the generation process as feedback. The exploitable parameters include gradient information (Zou et al., 2023; Li et al., 2024a), token probability distributions (Andriushchenko et al., 2025; Meng et al., 2025), and other internal representations (Labunets et al., 2025; Guo et al., 2024a). These white-box and gray-box approaches provide precise optimization signals by retrieving information related to the model’s internal states, and can therefore lead to an efficient optimization process.

The iterative approaches show better optimization potential, as the one-shot paradigm lacks feedback analysis. The capability of adaptation provided by feedback information determines that iterative attacks can adjust the prompts gradually closer to optimal and can thus achieve higher ASR. Furthermore, the multi-turn nature of iterative attacks creates opportunities for multi-turn-based jailbreak strategies like context priming and incremental boundary testing, etc. These strategies enhance the efficiency of iterative attacks. However, iterative approaches face various practical limitations. Firstly, iteratively accessing the model with a similar behavioral mode can trigger the defensive mechanisms. Moreover, both subcategories of iterative attacks suffer from insufficient data issues in practical scenarios: model response dependent approach suffer from vague directed searching due to limited feedback quality, thus often require large number of iterations to achieve success; model parameter

dependent approaches suffer from the strict data policy of the service providers, and often become invariable due to prohibit access to the required internal data.

**Parallel:** Parallel approaches involve the simultaneous deployment of multiple independent attack agents (Chao et al., 2024; Andriushchenko et al., 2025). Unlike iterative approaches, in parallel configurations, the search space is divided and assigned among the agents, so that the agents can explore different regions of the attack vector spaces, such as different prompt templates, varied optimization algorithms, complementary search heuristics to maximize coverage of potential vulnerabilities, or simply different prompt initializations.

This approach has the advantage of high time efficiency and attack coverage. Multiple parallel searches can discover effective jailbreak prompts faster than sequential approaches in ideal situations, since the diverse exploration strategies help agents avoid becoming stuck in local optima, which might trap individual agents but will not have a great influence on the whole process. Furthermore, parallel settings create a comprehensive vulnerability map through their extensive exploration of attack vectors against the target model.

However, parallel attacks face practical limitations and constraints. The most notable drawback is the explosive computational resource requirement. With multiple instances of attack algorithms and LLM completions running, the computational cost would be explosively large. While using APIs provided by the LLM providers can reduce local computational stress for completions, API rate limits and the high cost of API still cause problems in real-world scenarios.

Meanwhile, parallel configurations have low stealthiness, since simultaneous multiple queries from the same source may trigger the service provider’s security monitoring systems designed to identify coordinated attacks in real-world situations. Furthermore, coordinating multiple agents is also a challenge: Too many agents could lead to redundant exploration between agents and waste computational resources, whereas too few agents could result in insufficient attack coverage.

## C.2 Semantic

Both subcategories show distinct trade-offs in effectiveness and stealthiness.

**Natural:** Natural approach aims to jailbreak the target LLM with human-readable prompts.

Such prompts are fluent in expression and logic, which helps to circumvent detection mechanisms like PPL-based detectors. However, these attack prompts should still express the malicious intention, which requires effective camouflaging techniques (e.g., Role-playing, authority endorsement, etc.) (Shah et al., 2023; Jin et al., 2025; Tian et al., 2024; Bethany et al., 2024; Ding et al., 2024; Li et al., 2024c; Zeng et al., 2024). The core principle of these methods is replacing or extending the context of original harmful prompts with natural language to mask them as natural components of normal interactions.

The natural language form provides extensive sustainability to these methods. Yet, whether these methods are viable in real-world situations depends on their camouflaging techniques.

**Perturbed:** Perturbing approaches jailbreak target LLMs with ostensibly nonsensical sequences in the input prompts. To circumvent the model’s refusal without changing the malicious intention of the harmful prompts, two notable techniques have been applied:

Cipher-based methods employ various encoding techniques, including self-defined and well-known existing codings, to mask harmful content into a human-unreadable form (Yuan et al., 2024; Lv et al., 2024; Liu et al., 2024f; Jiang et al., 2024). With these techniques, attackers can avoid being detected by matching defenses that rely on keyword detection or content classification.

Another technique is suffix distraction. These methods usually employ optimization processes to filter out the optimal suffixes that can lead to positive model output according to the given harmful prompts (Andriushchenko et al., 2025; Zou et al., 2023; Li et al., 2024a). However, the optimized suffixes are usually garbled with no semantic meaning, reducing the stealthiness of the attack prompts.

Perturbing approaches enable more direct optimization and can achieve higher ASR through precisely targeting the model’s outputs. However, the reduced semantic meaning in the resulting prompts undermines the sustainability and also the feasibility of these methods in practical situations due to their low stealthiness and the rapid development of detection mechanisms.

### C.3 Model Accessibility

This dimension distinguishes the jailbreak methods according to the required accessibility level for the involved LLMs.

**Attacker Model Access:** Attacker model access refers to the accessibility of the models that are involved in the jailbreak pipelines, such as for jailbreak construction, optimization, and execution. Open-source attacker models sometimes face feasibility challenges, since local deployment may require the user to have the same instruments and configurations as the researchers. However, they provide better sustainability in the long run compared with closed-source models, as commercial licensing changes, service discontinuations, or policy modifications by service providers may influence closed-source model applications (Li et al., 2025; Liao and Sun, 2024; Paulus et al., 2025; Guo et al., 2024a; Liu et al., 2024c). Moreover, open-source models provide opportunities for extensive customization and fine-tuning, allowing for higher task adaptation.

Conversely, closed-source attacker models usually show better performance on complex tasks due to their commercial characteristics and associated resource investments (Chao et al., 2024; Li et al., 2024c; Kuo et al., 2025). The commercial nature of these models brings them better computational resources, specialized research teams, dedicated datasets, and effective training methods. All these can be economically unfeasible for open-source model developments. On one hand, only large commercial companies can afford the expensive equipment and large-scale training required to support much larger models; on the other hand, commercial incentives motivate continuous investment in equipment and cutting-edge technical development to maintain market competitiveness. For these reasons, closed-source models typically maintain higher reasoning abilities, context understanding, and instruction following abilities that are essential for complex jailbreak prompt generation. These abilities lead to higher performance in nuanced manipulation (e.g., tone, attitude, etc.) of the prompts, making the attack more effective. Although the easy deployment of closed-source models via API provides high feasibility, their application is accompanied by multiple operational risks. Access restrictions, pricing changes, usage policy modifications, or complete service termination could all severely harm the methods’ long-term sustainability. Additionally, closed-source models limit transparency in the attack prompt generation phase. With less evidence to explain the attack success, methods relying on closed-source attacker models are often considered less robust than open-source ones, and

1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921

this could hinder methodological improvements and scientific understanding.

The choice of attacker model access presents a trade-off between feasibility versus sustainability and attack complexity versus customization.

**Target Model Access:** Target model access refers to the access level to the target model required by the attack methods. Access level can influence the available attack strategies greatly. Open-source target models provide comprehensive model internal information, such as architectural details, parameter weights, training data characteristics, and intermediate activation patterns (Zou et al., 2023; Liao and Sun, 2024; Li et al., 2024a; Guo et al., 2024a). This data then enables white-box attacks that can exploit specific architectural vulnerabilities, manipulate internal representations, perform gradient-based optimizations, and develop targeted perturbations for jailbreak. In short, the transparency of open-source models provides much jailbreak-relevant data and can greatly help the attacker to manipulate the model’s response in various ways, improving the sustainability.

In contrast, closed-source target models restrict attackers to black-box interaction with API endpoints and gray-box interactions in user interfaces (Wu et al., 2025; Zeng et al., 2024; Shah et al., 2023). In this situation, only very limited or even no internal data of the models is available to the attackers, making the design of the jailbreak attack strategies harder. Moreover, in real-world situations, commercial closed-source target models are usually equipped with extra comprehensive defensive structures, such as proprietary safety filtering systems, advanced content moderation mechanisms, behavioral monitoring algorithms, and dynamically updated safety policies, that are also inaccessible to external attackers. Due to the commercial deployment context, where in many situations failure is not acceptable, commercial model providers invest heavily in safety systems to protect their market position, maintain user trust, comply with regulatory requirements, and avoid potential legal problems. Circumventing these complicated and nested defensive mechanisms indicates a higher feasibility in the real world compared to circumventing the relatively simple and transparent defensive mechanisms of the open-source models.

Moreover, successful jailbreaks against closed-source models have greater practical meaning for real-world security assessment. Since these models are deployed in commercial environments with

real users, successful attacks can point out vulnerabilities in systems that actually matter for practical security, rather than purely academic research scenarios. This makes jailbreak success on closed-source target models more valuable for both offensive security research and defensive mechanism improvement.

The accessibility dimension creates a strategic trade-off matrix where different combinations of attacker and target model access levels result in distinct attack paradigms.

Open-source attacker models targeting open-source systems enable the most comprehensive and sustainable attacks, but may have limited real-world feasibility due to the deployment difficulty and ideality of the target.

Open-source attacker models targeting closed-source systems offer a balance of sustainability and practical relevance, allowing researchers to develop sustainable methods while demonstrating vulnerabilities in commercially deployed systems.

Closed-source attacker models targeting open-source systems leverage superior capabilities for attack generation while benefiting from transparent target analysis, though this combination may have limited practical significance.

Closed-source attacker models targeting closed-source systems represent the most realistic and high-stakes attack scenarios, combining superior attack capabilities with valuable target systems, showing high feasibility. But it suffers from significant sustainability challenges due to dependencies on commercial services for both attack generation and target access.

#### C.4 LLM Role

This Dimension distinguishes methods according to what role LLMs have played in the attack pipeline. The role classification influences attack methodology, robustness, and adaptation.

When LLMs are deployed as **generators**, they are responsible for jailbreak prompt creation, modification, and optimization. Generator-based approaches rely on the strong linguistic ability and reasoning capabilities of LLMs to exploit the vulnerabilities of target LLMs and circumvent defensive mechanisms. To be more specific, generator LLMs are good at producing human-like, contextually coherent attack prompts that are less likely to trigger simple defenses such as keyword matching, and can automatically generate attack prompts at a large scale without manual prompt engineering.

1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972

1973 However, generator-based systems face multiple  
1974 limitations. Firstly, regardless of whether they are  
1975 open-source or closed-source models, the safety  
1976 alignment update, accompanied by regular model  
1977 updates, tends to reduce the possibilities of the  
1978 model generating harmful content. This means that  
1979 the generator models will gradually refuse to gener-  
1980 ate attack prompts that contain malicious intentions  
1981 or try to avoid generating harmful content in their  
1982 response. This reduces the sustainability of these  
1983 generator-reliant methods. Moreover, generator-  
1984 based methods face generation quality issues. To  
1985 explore larger attack vector spaces, LLM temper-  
1986 atures are typically set to non-zero values, which  
1987 can lead to stochastic and unpredictable output.  
1988 This results in variation in the effectiveness of the  
1989 attack, and could undermine the real-world feasi-  
1990 bility of these methods; besides, the coverage of  
1991 the training data of the generator model may be  
1992 imbalanced in some specific domains. This results  
1993 in performance variation in different scenarios and  
1994 benchmarks. Finally, generator approaches face  
1995 increased exposure risks, as the generated attack  
1996 prompts may exhibit characteristic patterns of the  
1997 generator LLMs, such as high-frequency words and  
1998 tones, specific logical patterns, etc. These can be  
1999 identified by gradually developing defensive mech-  
2000 anisms, thus would reduce the sustainability of the  
2001 methods.

2002 When LLMs function as **evaluators**, they work  
2003 as the assessment mechanism in the attack pipeline.  
2004 They analyze the target model’s response and the  
2005 generated attack prompts and provide information  
2006 to guide the prompt optimization process. Through  
2007 this analysis, the evaluator models can determine  
2008 whether there is policy-violating content in the  
2009 model’s response, or harmful content generation  
2010 evasions, etc. It can then identify partially success-  
2011 ful attempts that can be refined, and direct search  
2012 algorithms toward more promising attack vectors.  
2013 In short, evaluator LLMs’ strong analysis capabil-  
2014 ities can be used to reduce the search spaces of  
2015 optimal attack vectors and accelerate the attack  
2016 prompt optimization phase.

2017 A key advantage of applying evaluator LLMs is  
2018 that they can identify differences between target  
2019 models’ declared safety policies and their actual  
2020 implementation. Evaluator models help attackers  
2021 to focus on the characteristics of harmful prompts  
2022 that can lead to harmful outputs, i.e., those attack  
2023 vectors that may more easily achieve a jailbreak.  
2024 This helps the attackers to figure out the precise

2025 vulnerabilities of the target models and make di-  
2026 rected improvements. On the other hand, this is  
2027 also important for defensive mechanism develop-  
2028 ment, since evaluator LLMs provide information  
2029 regarding the weakness of the target model, which  
2030 should be complemented by the defensive mecha-  
2031 nisms. Also, unlike the generator models, evaluator  
2032 models are less sensitive to model updates since  
2033 they are only used for text analysis. The researchers  
2034 only need to update the rule embedded in the in-  
2035 struction prompt or update the finetuning dataset  
2036 to keep the evaluator model up to date, thus giving  
2037 the methods great sustainability.

2038 However, the application of evaluator models  
2039 in the jailbreak attack pipeline also faces various  
2040 challenges. The first challenge is that the evaluator  
2041 model may misclassify jailbreak successes. Since  
2042 the feedback from the target model does not di-  
2043 rectly provide the information that leads to optimal  
2044 attack prompts, it could cause the evaluator model  
2045 to over- or underestimate some attack vectors, thus  
2046 preventing the attack prompts from being optimal  
2047 and limiting the effectiveness of the attack meth-  
2048 ods. Moreover, evaluator models also face the same  
2049 data imbalance challenge as the generator model,  
2050 limiting the generalization of a single evaluator  
2051 model in different scenarios. Also, deploying eval-  
2052 uator models means extra computational resource  
2053 requirements, which reduces their real-world feasi-  
2054 bility when the attempt number is very high.

2055 In summary, the generator offers attack sophis-  
2056 tication and creativity. Generator LLMs can ex-  
2057 plore novel attack vectors and adapt to specific  
2058 target characteristics. However, generators suffer  
2059 from safety alignment updating, inconsistent output  
2060 quality, computational overhead, and increased de-  
2061 tectability compared to manually crafted prompts.  
2062 Conversely, the evaluator role provides stability  
2063 and maintainability. The evaluator helps to filter  
2064 out effective attack vectors. However, evaluators  
2065 face accuracy challenges, domain blind spots, and  
2066 computational overhead.

### D Paper Distribution

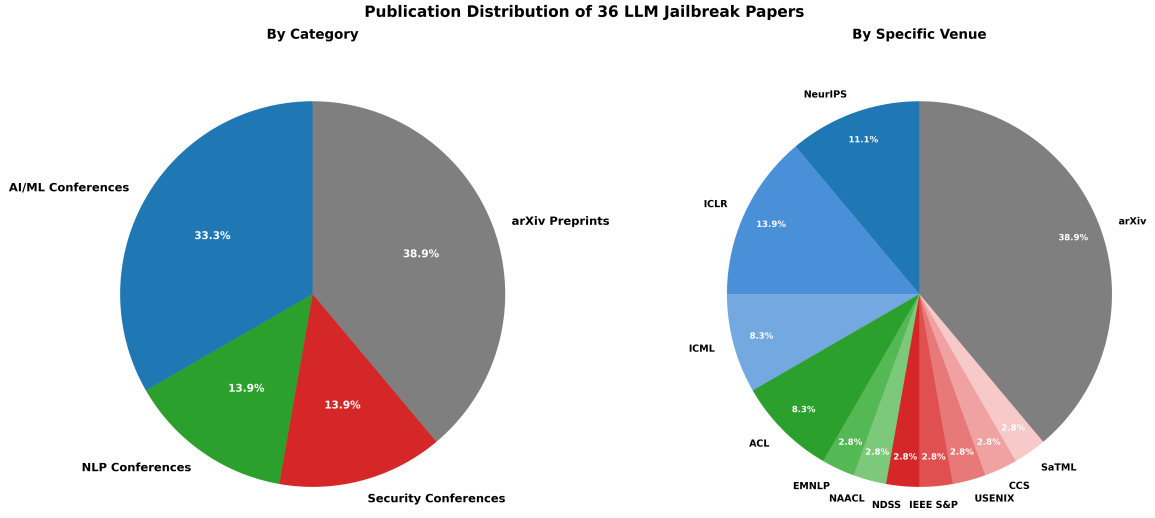


Figure 3: Venue Distribution of Surveyed Papers

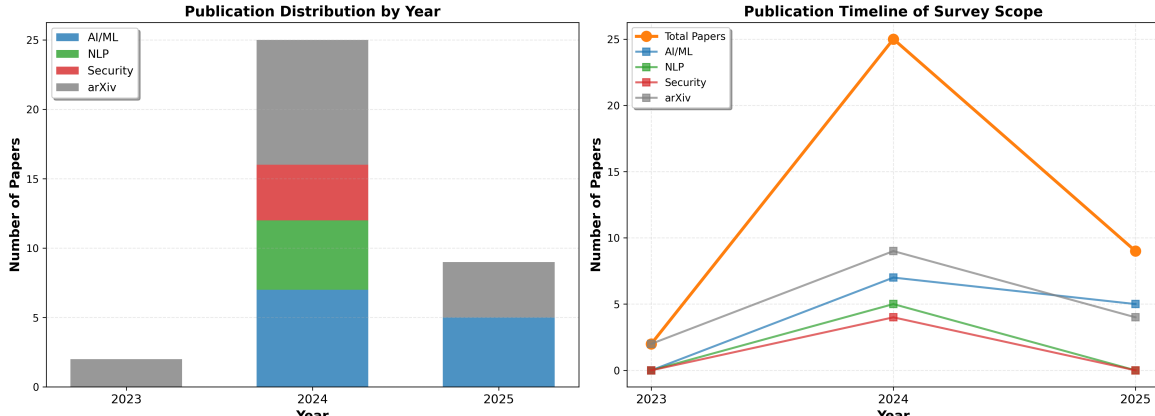


Figure 4: Time Distribution of Surveyed Papers