

# An Empirical Study on Prompt Injection Attacks and Defences

Haritha Gunarathna

*Department of Computer Engineering*  
*University of Peradeniya*  
Kandy, Sri Lanka  
e18118@eng.pdn.ac.lk

Denuwan Weeraratne

*Department of Computer Engineering*  
*University of Peradeniya*  
Kandy, Sri Lanka  
e18382@eng.pdn.ac.lk

Nimuthu Wijerathne

*Department of Computer Engineering*  
*University of Peradeniya*  
Kandy, Sri Lanka  
e18398@eng.pdn.ac.lk

Asitha Bandaranayake

*Department of Computer Engineering*  
*University of Peradeniya*  
Kandy, Sri Lanka  
asithab@eng.pdn.ac.lk

Roshan Ragel

*Department of Computer Engineering*  
*University of Peradeniya*  
Kandy, Sri Lanka  
roshanr@eng.pdn.ac.lk

Damayanthi Herath

*Department of Computer Engineering*  
*University of Peradeniya*  
Kandy, Sri Lanka  
damayanthiherath@eng.pdn.ac.lk

**Abstract**—The recent surge of Large Language Models (LLMs) has revolutionized Natural Language Processing (NLP), showcasing remarkable abilities in text generation, translation, and comprehension. However, many vulnerabilities have been detected in LLMs. Prompt injection attacks specifically pose a significant threat to the security and integrity of LLMs. This paper presents a comprehensive empirical study on this subject. We analyze existing research on prompt injection attack models, exploring their nature, attack models and prevailing taxonomies. We also delve into the impact of these attacks and discuss the current/proposed defense strategies and mitigation. Finally, we discuss potential future directions in the area of prompt injection attacks. By this paper, we hope to contribute to the ongoing effort of securing LLMs and advancing responsible AI development.

**Index Terms**—Large Language Models (LLMs), Prompt Injection Attacks, Adversarial Attacks, Security in LLMs, Jailbreaking, Indirect Prompt Injection

## I. INTRODUCTION

With its remarkable ability to produce text that is both coherent and contextually relevant, Large Language Models (LLMs) have brought in a new era of natural language processing. This advancement hasn't, however, happened without scrutiny, and one major concern is that LLMs are vulnerable to adversarial attacks [34], [39]. Among these, prompt injection attacks have emerged as a particularly complex and potent form of manipulation.

Prompt injection attacks involve the strategic crafting of input prompts to deceive LLMs into producing unintended, unethical or biased outputs. They take use of the adaptability of language models that depend on user-provided prompts to generate responses, in contrast to classic adversarial attacks, which concentrate on altering specific tokens or words. This raises critical questions about the robustness and security of LLMs, especially as they are diverse in their applications, ranging from chat bots to content generation.

This review delves into the intricacies of prompt injection attacks within the broader context of adversarial attacks

on LLMs. By examining the attack methods, impact, classifications and potential mitigation prevailing in the existing literature, this exploration aims to shed light on the evolving landscape of LLM security focusing on prompt injection attacks and contribute to the ongoing dialogue surrounding responsible AI development and deployment.

## II. METHODOLOGY

In this section, we outline the methodology employed for the collection, testing, and documentation of LLM security and various prompt injection attacks and defenses. Given the primary objective of providing a comprehensive overview of the existing literature, our approach involved a systematic review of diverse sources, including arXiv preprints and platforms such as jailbreakchat.com, learnprompting.org, and owasp.com. Given the dynamic nature of the field and the frequent updates to chatbots and LLM APIs, we placed substantial emphasis on non-academic sources.

To identify prompt injections, we conducted keyword searches on Google, Google Scholar, and arXiv, employing terms such as “prompt injection”, “jailbreak”, “large language models (LLMs)”, “direct prompt injection attacks”, “indirect prompt injection attacks”, and “LLM defense”. Complementing academic sources, we thoroughly reviewed content available on the internet, including websites, social media and blogs. Addressing concerns related to the lack of peer review in non-academic sources as this is a relatively new subject, we adopted a two-step verification process. Initially, we sought multiple sources making the same claim, supported by credible screenshots illustrating the prompt injection mechanism. Subsequently, we conducted prompt injection tests in Bard, Bing Chat, GPT-3, and GPT-4 models.

Due to the swift patching of vulnerabilities, these methods do not consistently work without additional prompt engineering or may no longer be effective at the time of writing. As such,

these vulnerabilities are pertinent primarily to future designs of chatbot and LLM interfaces.

### III. REVIEW FINDINGS

#### A. LLMs, advancement, and usage

With their exceptional ability to comprehend and produce human language, Large Language Models (LLMs) have taken their place as a key instrument in the Natural Language Processing (NLP) domain. Mathematical concepts based on deep learning architectures, like transformer models or recurrent neural networks (RNNs), are at the core of LLMs. LLMs can be categorized into various types, including transformer-based models (e.g. GPT-3), recurrent neural networks (e.g. LSTM), or even models using novel architectures such as the Transformer-XL, as shown in [18]. They have been extensively deployed across numerous fields and will be essential components in future communication networks [18]. Models (LLMs) like GPT-4 [6], LLaMA [11], and Bard [1], have dramatically transformed a wide array of applications with their exceptional ability to generate human-like texts [12]. Their applications span a broad spectrum of tasks such as machine translation, sentiment analysis, question answering, and text summarizing, opening new avenues for innovation and research in the field [18]. Large Language Models (LLMs) have been revolutionizing our lives in many ways, not just for practitioners and scholars but also for the general public and will continue to do so. This is evident by ChatGPT, which shortly after its release, gained immense popularity, attracting over 100 million users in a short period of time [7]. Furthermore, there is a constant stream of new models, including the more advanced GPT-4 [6], being introduced at a quick pace, and smaller white-box models [13].

Due to their superb generative capability, LLMs are widely deployed as the backend for various real-world applications called LLM-Integrated Applications. For example, Microsoft utilizes GPT-4 as the service backend for new Bing Search [1]; OpenAI developed various applications—such as ChatWithPDF and AskTheCode—that utilize GPT-4 for different tasks such as text processing, code interpreter, and product recommendation [16]. Other third parties can utilize LLMs for their diverse applications as well, including chatbots, assistants, and various other applications. With the establishment of the GPT Store in ChatGPT, custom models that align with the person's or organization's unique requirements and data have become publicly accessible, creating a marketplace of various AI tools designed for various applications [28].

#### B. Vulnerability Analysis of Generative AI models

The quick rise of LLMs and their expanded usage have given rise to many security concerns and vulnerabilities. In adversarial attacks [15], [26], which is a known threat to machine learning algorithms, carefully manipulated inputs can drive a machine learning structure to produce reliably

erroneous outputs to an attacker's advantage [19]. Attacks can be targeted, seeking to change the output of the model to a specific class or text string, or untargeted, seeking only to result in an erroneous classification or generation as shown in [19]. Jailbreak, backdoor attacks, and complex data poisoning is shown as some examples of the broad spectrum of adversarial attacks in LLMs as per [12].

In [19] adversarial attacks are classified according to 5 concepts, Learning Structures, Injection source, Attacker access, Attacker type and the attack goal.

Based on the learning structure, attacks can be classified into uni-modal, multi-modal and additional attacks. Jailbreak attacks and prompt injection attacks are the two prevalent types of adversarial attacks on aligned uni-modal Large Language Models, where uni-modal refers to focusing on manipulating a single type of data or output modality [19]. Multi-modal attacks, which are models that accept as input not only text, but additional modalities such as audio or images, are divided into 4, Manual attacks, systematic adversarial attacks, white-box attacks and black-box attacks. Other prevalent attacks can be classified under Additional Attacks as per the learning structure, which can be further divided into adversarial attacks in complex systems and earlier NLP attacks.

#### C. Prompt Injection Attacks and their Impact

Prompt injection attacks refer to a type of security threat where malicious users manipulate the prompts provided to Large Language Models (LLMs) or other AI systems to influence the generated outputs in unintended ways [1], [35]. This method involves crafting input prompts in a manner that bypasses the model's safeguards or triggers undesirable outputs. It is the OWASP No.1 threat for LLMs [8]. There have been many adversarial effects caused due to prompt injection in generative AI models, some of which are listed below.

#### **Generating text classified under prohibited scenarios**

8 distinct prohibited scenarios from OpenAI's disallowed usage policy [6] that represents potential risks and concerns associated with the use of ChatGPT, was attempted to jailbreak in [22]. Comparatively [30] created and deployed a question set comprising 46,800 samples across 13 forbidden scenarios. Some forbidden scenarios taken into experimentation were Harmful content, Illegal activities, Political Lobbying, and Adult content. Both experiments showed that current LLMs and safeguards cannot adequately defend jailbreak prompts in those scenarios. Among the 13 forbidden scenarios tested in [30], Political Lobbying was found to be the most vulnerable to jail-breaking, followed by Pornography and Legal Opinion. With the prompts used in [22] it was found that Illegal activities, Fraudulent or Deceptive activity and Adult content were the easiest scenarios to be broken by jailbreak prompts. GPT-4 demonstrated a greater resistance against jailbreak prompts aimed at extracting prohibited content, compared to GPT-3.5-TURBO [22]. Additionally, it was concluded that jailbreak prompts easily achieve high success rates even in

scenarios where initial resistance is observed, proving the capability of jail-breaking [30]. This was shown in [22] as well, as jailbreak prompts achieved a success rate of 74.6%, which is much higher than the rate of 29.0% for non jailbreak prompts.

**Goal Hijacking** Goal Hijacking, also known as “Prompt Divergence” [36] attempts to redirect the LLM’s original objective towards a new goal desired by the attacker [19]. In other words, malicious users input prompts that hijack the original goal of LLM-integrated applications [17]. By deploying a novel framework ‘PROMPTINJECT’ in [14] to assemble prompts in the aim of goal hijacking, which achieved a success rate of 58.6%, it is shown that a malicious user can easily perform goal hijacking via human-crafted prompt injection. [19] Inspired by data poisoning and backdoor attacks, [37] introduce a novel concept of “Virtual” prompt injection attacks, focusing on goal hijacking, causing the model to answer a different question resulting in an answer of use to the attacker. Remarkably it was shown that, by contaminating only a small fraction of the instruction-tuning dataset, the attacker can influence the model’s behavior during inference when the model is queried about a specific target topic.

**Prompt leaking/ System prompt extraction** [14] defines prompt leaking as the act of misaligning the original goal of a prompt to a new goal of printing part of or the whole original prompt instead. A malicious user can try to perform prompt leaking with the goal of copying the prompt for a specific application, which could be the most important part of GPT-3-based applications. If an attacker can get access to the system prompt of a service provided by a company, they can build a clone of the service using the recovered system prompt, making this prompt a valuable part of each system’s intellectual property [19]. Demonstrating this, [29] employs a leaked prompt to construct a mock LLM integrated application, observing a high degree of functional similarity between the two applications, implying that the leaked prompt can effectively replicate the capabilities of the original application as said before. A success rate of 23.6% was achieved through ‘PROMPTINJECT’ for prompt leaking [14], showing that it is notably more challenging than goal hijacking. However, [28] observed that deploying prompt injection attacks on 200+ custom GPTs yielded an alarming 97.2% success rate for system prompt extraction. They have defined system prompt extraction as the act of deceiving custom GPTs into disclosing the designed system prompt, which is the same as prompt leaking.

**File leakage** File leakage is the act of stealing the designer-uploaded files used by the custom GPT for its specific purpose [28]. It has been identified that this vulnerability enables a malicious user to detect files uploaded by the custom GPT developer, including identifying the names and sizes of these files. The attacks on 200 custom GPTs, aiming file leakage,

has a success rate of 100%, highlighting an immediate need to address the risk of prompt injection.

**Hallucination** LLMs might make up facts (“hallucinate”), generate polarized content, or reproduce biases, hate speech, or stereo types [22], partially stemming from pre-training on massive crawled data sets. Even though one of the motivations for leveraging Reinforcement Learning from Human Feedback (RLHF) is to better align LLMs with human values and avert these unwanted behaviors [38], OpenAI reports that GPT-4 shows a tendency, to still hallucinate or generate harmful content [39].

**Social Engineering/ Information gathering** Attackers can manipulate the “instructions” given to generative AI models, tricking unsuspecting individuals to reveal personal information, clicking malicious links, or sharing misinformation further, causing potential financial loss, reputation damage, or societal disruption [9]. As shown in [13] Indirect prompt injection could be leveraged to exfiltrate users’ data (e.g., credentials, personal information) or leak users’ chat sessions. This can be done in interactive chat sessions by persuading users to disclose their data or indirectly via side channels. ChatGPT’s ability to understand context, impressive fluency, and mimic human-like text generation could be leveraged by malicious actors [30]. The power of this approach lies in ChatGPT’s ability to generate text that aligns with the victim’s expectations, thereby increasing the likelihood of the victim complying with the request. [30]. Exfiltrating user’s private information is also listed as a threat by indirect prompt injection attacks in [17].

#### *D. Characterization and Taxonomy of Prompt Injection Attacks*

Current prompt injection attacks predominantly fall into two categories [35], direct and indirect prompt injection. In the case of direct prompt injection, an attacker has the capability to control their own engagement with the LLM. On the other hand, with indirect (or cross-domain) prompt injection, the attacker possesses the ability to manipulate the interaction of another user [33]. This taxonomy and further classifications are shown in Figure 1.

**Direct prompt injection attacks** Direct prompt injection attacks [2], [35] operate on the premise of a malicious user directly injecting harmful prompts into the application inputs. The primary goal of these attacks is to manipulate the application into responding to a different query than its original intent. To achieve this, the adversary creates prompts, which can be natural language or communicate with LLM using cipher codes [25] like Caesar Cipher or Binary that can influence or negate the predefined prompts in the combined version, leading to the desired responses.

We found that there exists a debate on whether jailbreak attacks [22] fall under direct prompt injection attacks as mentioned in [35]. In some papers [33], [35], jailbreak attacks

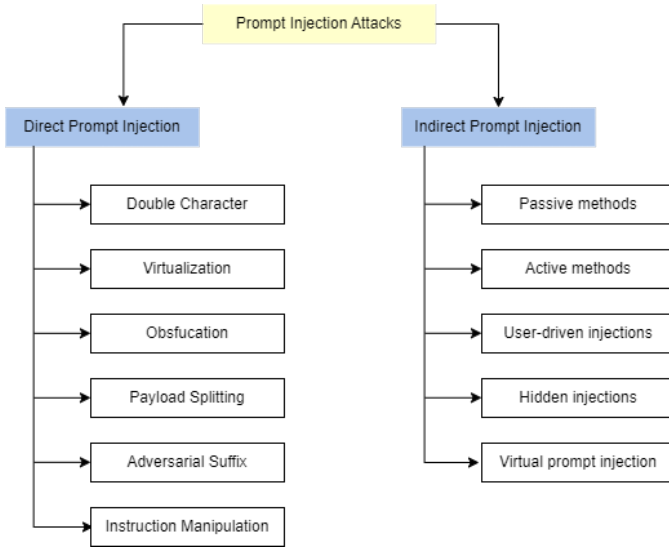


Fig. 1. Taxonomy of existing prompt injection attacks

are not considered as a prompt injection attack. Difference they highlight is that jailbreak attacks typically involve eliciting model-generated content that divulges training data specifics, potentially leading to privacy breaches. While direct prompt injection attacks focus on manipulating the LLM's output by appending malicious prompts to the input sequence, jailbreak attacks [3] target the model's training data specifics to extract sensitive information. But when referring to web pages like owasp.org [36] it clearly mentioned that jailbreak attacks are a type of direct prompt injection attacks.

According to [36] there are six sub classes of direct prompt injection attacks. They are Double Character, Virtualization, Obfuscation, Payload Splitting, Adversarial Suffix and Instruction Manipulation.

These sub classes are described in Table 01. In addition, further examples and experiments done across different LLMs are displayed in Table 03 (Appendix A).

**Jail-breaking** Further research has been done on Jail-breaking as a separate topic. 10 patterns of Jail-breaking have been identified in [22], spread across 3 categories, Pretending, Attention Shifting and Privilege escalation. It is evident from this study that pretending is the most prevalent strategy used by attackers to bypass restrictions while Attention shifting and Privilege escalation are less frequently employed. Jailbreak prompts have been classified to 8 communities in [30], each demonstrating diverse and creative attack attempts in designing jailbreak prompts. It's found that Jailbreak prompts have evolved to be more stealthy and effective to conceal their malicious intent, evident by their reduced length, increased toxicity and semantic shift [30]. GPT4 demonstrates a greater resistance against jailbreak prompts aimed at extracting prohibited content, compared to GPT3.5 turbo, mainly because GPT4 has an improved ability to comprehend the output meaning [22]. Even though LLMs

Type of Direct PI Attack	Description
Double Character	A prompt that produces a response consisting of one character constrained by language model rules and another character unrestricted, bypassing content limitations.
Virtualization	Create a scenario where the LLM operates in an unrestricted mode, resembling a developer environment or a virtual scenario where potentially harmful content is generated within a "virtual machine."
Obfuscation	Create a prompt where malicious or rule-breaking instructions are obscured, such as encoding them as base64 characters instead of regular ASCII characters.
Payload Splitting	Combine benign instructions from multiple prompts into a final prompt, where individually harmless texts become malicious when merged together, like A and B.
Adversarial Suffix	Generate a random suffix composed of words and characters to append to a malicious prompt, effectively bypassing alignment mechanisms of the LLM and eliciting a response to the malicious prompt.
Instruction Manipulation	Craft a prompt that either discloses the pre-written instructions or the original prompt provided to the LLM interface, or instructs the interface to disregard said instructions.

TABLE I  
DIRECT PROMPT INJECTION ATTACKS

trained with RLHF (Reinforcement Learning based on Human Feedback) show initial resistance to forbidden questions, they have weak resistance towards jailbreak prompts [30] External safeguards, such as OpenAI Moderation Endpoint and NeMo-Guardrails, demonstrate limited capability of identifying jailbreak prompts hence calling for stronger and more adaptive defense mechanisms [30].

**Indirect prompt injection attacks** The concept of Indirect Prompt Injection to compromise LLM-integrated applications was introduced in [13]. The first taxonomy and systematic analysis of the threat landscape associated with Indirect prompt injections in LLM-integrated applications was developed in [13] as well. In indirect prompt injection attacks, attackers inject malicious instructions into third party content, which when retrieved by an LLM-integrated application and ingested by the LLM, cause the LLM's output to deviate from the user's expectations [17]. As described in [13], a challenge associated with Large Language Models (LLMs) arises from the incorporation of retrieval, introducing ambiguity between data and instructions. Consequently, adversaries gain the capability to remotely influence other users' systems by strategically embedding prompts within data expected to be retrieved during inference. The LLM exhibits a reduced capacity to discern concealed instructions within the data, as the injection of prompts into the data becomes a viable operation. Hence it leads to Indirect prompt Injection attacks

Type of Indirect PI Attack	Description
Passive Methods	Strategically embed malicious prompts or content within a publicly accessible source, such as webpages, or social media posts which could be scanned and interpreted by LLMs. This encompasses the manipulation of data, particularly text, that is evaluated by LLMs when browsing the web.
Active Methods	Deliver malicious prompts directly to an LLM, for example by sending emails containing such prompts, exploiting email clients equipped with LLM extensions to execute the prompt upon receipt.
User-driven Injections	Utilize social engineering techniques to share seemingly innocent prompts, persuading unsuspecting users to copy and paste them into an LLM, unaware of the underlying malicious intent.
Hidden Injections	Attackers may employ multiple exploit stages. Initially, a smaller injection instructs the model to fetch a larger payload from another source.
Virtual Injections	The attacker manipulates the instruction tuning data of an LLM, causing specific scenarios where the model's behavior becomes misaligned, resulting in outputs as if it received additional instructions through a prompt.

TABLE II  
INDIRECT PROMPT INJECTION ATTACKS

[13]. As LLMs progress, they demonstrate the capability to reinforce their response to a user prompt by surfing popular websites like Wikipedia etc. [13] explains that it is done in a self-supervised manner employing in-context learning or using ReAct, utilizing Chain-of-Thought prompting. Hence with the progress of Large Language Models (LLMs) and their increasing capabilities, there is a growing susceptibility to indirect prompt injection attacks. Specially, LLM-integrated applications are vulnerable to indirect prompt injection attacks where malicious LLM misbehavior response instructions embedded within external content compromise LLM's output, causing their responses to deviate from user expectations [17]. There have been 5 methods for indirect prompt injection attacks introduced in [13] and publicly accepted, being passive injections, active injections, user-driven injections, and hidden injections. Another indirect prompt injection attack method, as described in [42] is Virtual injection attacks. These categories are described in detail in Table 02. Indirect prompt injections have resulted in many adverse threats like Information gathering, fraud and spreading malware [13].

#### E. Novel Attack models

Many prompt injection attack models have been designed and introduced in research papers, some being entirely novel attacks that draw inspiration from the general context [14], [29], while others are frameworks that try to systematize and formalize prompt injection attacks [16].

A novel prompt injection framework called 'PROMPTINJECT' was proposed in [14], for mask-based iterative adversarial prompt composition and it was shown how GPT3, the most

widely deployed language model in production, can be easily misaligned by simple handcrafted inputs. In particular, two types of attacks are demonstrated, goal hijacking and prompt leaking. The framework consists of 2 parts, the base prompt and the attack prompt.

A novel black box prompt injection attack technique called HOUYI was introduced in [29], drawing inspiration from traditional web injection attacks. HOUYI is compartmentalized into three crucial elements: a seamlessly-incorporated pre-constructed prompt, an injection prompt inducing context partition, and a malicious payload designed to fulfill the attack objectives. Leveraging HOUYI, previously unknown and severe attack outcomes were unveiled, such as unrestricted arbitrary LLM usage and uncomplicated application prompt theft. HOUYI was deployed on 36 actual LLM-integrated applications, where it was successful on 31 applications. HOUYI consists of a prompt refinement with dynamic feedback, ultimately, helping it to output a collection of successful attack prompts.

In [43], an innovative framework called MASTERKEY, that reverse engineers the hidden safety mechanisms in LLMs against jailbreaking and is capable of generating jailbreak prompts that work across multiple mainstream LLM chatbots (including CHATGPT, Bard, Bing Chat) is introduced. MASTERKEY begins by disassembling the jailbreak defense mechanisms used by different LLM chatbot services, paying particular attention to the relationship between the length of the LLM's response and its generation time. A time-based LLM testing strategy is developed by using this correlation as an indicator and utilizing the blind SQL attack mechanism found in traditional web application attacks. A unique methodology to automatically generate universal jailbreak prompts is introduced, based on these features and findings, achieving an average success rate of 21.58% among mainstream chatbot services. An important finding in this study is that the existing jailbreak prompts seems to be effective towards CHATGPT only, while demonstrating limited success with Bing Chat and Bard. However the proposed attack framework shows success in those models as well.

As a different approach (and the first of its kind) a framework was proposed in [16], to systematize and formalize prompt injection attacks. Under the framework, different prompt injection attacks essentially use different strategies to craft the compromised data prompt based on the clean data prompt, injected instruction of the injected task, and the injected data of the injected task. Existing attacks are considered to be special cases in this framework. Based on this framework, a new prompt injection attack was designed by combining existing attack strategies, and the results show that this combined attack outperforms others. A defense framework was also introduced along with. The introduced framework-inspired attack was proved to be consistently effective for different target and injected tasks. Another interesting finding of this study [16] was that in general, the attack is more effective when the LLM is larger (or more powerful), and the reasoning that a larger LLM is more powerful in following the

instructions and thus is more vulnerable to prompt injection attacks, is suspected.

#### *F. Defense Mechanisms*

In the conducted extensive review over the available literature on the domain of Prompt Injection Attacks, it was found that there have been many attempts to come up with a relatively accurate and unbreakable defense mechanism against these types of attacks. Following is a brief description of the findings.

In paper [17], authors propose different types of Defense mechanisms based on two different aspects. One is White-box defense. This mechanism leverages the access to work with model parameters and architecture along with the training process. Other mechanism is called Black-box Defense, where there is no access to the model inner workings, parameters and training process, but it is based on prompt learning techniques such as usage of border strings, data-marking, usage of multi-turn dialogue and in-context learning (few-shot learning). White-box Defenses mainly includes fine tuning the model through adversarial training. Main motivation of these defense mechanisms is the LLMs' inability to distinguish between external content and instructions. [17]

As another defense mechanism, another paper [4], [17] propose a harmful filter be used in front of the inference model. The harm filter can be another LLM or another instance of the same LLM. The harm filter LLM is primed to identify the toxicity of the incoming prompt. The user input is formatted into a specific format before feeding into the harm filter. This method has shown significant evaluation metrics when tested with prominent LLMs like GPT3 and Llama2.

It is also possible to use more traditional ways of defending like input validation and output validation. Regex patterns, whitelists, blacklists can be used to identify a potential prompt injection attack [4], [5].

Post-prompting can also significantly improve the defense layer strength of LLMs. What this does is passing the system/user instruction after passing the external (potentially malicious) content first into the LLMs. This way it can work against attacks which convince the model to ignore all previous instructions [14].

Some other notable and traditional defense mechanisms are using XML tags (use XML tags to cover the external content), Random sequence enclosure (enclosing the external content with Random sequences) and sandwich defense (enclosure of the external content between the two instances of the same original prompt). All these methods try to enhance the LLMs ability to identify the line between provided instructions and the external content [1].

A different approach was taken by the author Xuchen [41] in his paper where he signs the user instruction prompt in a certain way that even the same prompt, when exposed to the model as two prompts (signed and unsigned), the model is able to identify as two different entities. This way, unless the signing signature is released to the public, a properly trained

LLM is capable to differentiate between malicious instructions and user instructions and act accordingly.

Other than this there were many other recently introduced defense mechanisms as well in addition to the traditional defense mechanisms we found when reviewing many other papers. [23], [27], [31], [32], [34], [40].

### IV. GAPS IN RESEARCH

Our comprehensive review of existing research on Prompt injection attacks and defenses has highlighted opportunities and areas for further exploration and research. They are given below.

#### *A. Lack of defense mechanisms for Transferable Direct prompt injection attacks*

Transferable Direct Prompt Injection is a more sophisticated version of prompt injection attacks where attackers study an LLM's responses to carefully crafted prompts, even seemingly non harmful ones. This analysis allows them to create triggering prompts that can reliably exploit LLM vulnerabilities even across different models. Developing techniques to detect and mitigate transferable attacks poses a significant challenge, as existing defense mechanisms designed for specific prompt injection methods may not be effective against transferable variants.

#### *B. Application-Specific Considerations*

As said earlier, LLMs are deployed as the backend of various applications for a wide range of purposes. The effectiveness of both attacks and defenses can vary significantly depending on the specific application and LLM used. A vulnerability analysis of different LLM architectures needs to be done, to understand how different LLM architectures and the applications built on top of them are susceptible to various prompt injection attacks, in order to develop targeted defense strategies. Developing defenses that can adapt to the specific context and task of the LLM application would improve their effectiveness. A further area needing research is the lack of an efficient method to separate User Prompt and information gained by external sources, which is a main method used in LLM-integrated applications.

#### *C. Lack of Systematic Understanding*

A lot of the current research primarily focuses on case studies and specific attack vectors. A systematic understanding of prompt injection attacks is lacking, hindering the development of comprehensive defenses. This can include classifying different attack types based on their goals, techniques, and impact which would provide a clearer picture of the threat landscape. Developing a formal framework for representing and analyzing prompt injection attacks would allow for more rigorous research and easier comparison of defense mechanisms.

## V. CONCLUSION

With the recent rise of Large Language Models or LLMs, they have made their way into the society, directly and as the back end of production grade applications and user interfaces. With this different adversarial attacks have emerged against LLMs along with various defence strategies to prevent and mitigate those.

In this study we have reviewed numerous peer-reviewed and preprint papers to map the existing literature into a clear and comprehensive introduction of Prompt Injection Attacks and Defenses. We start the review with an introduction to our empirical study, and then briefly describe the methodology we followed to search for existing literature on the topic. An important factor to consider is that most of the academic research papers we reviewed were not peer reviewed (preprints), mainly due to the fact that the topic of LLMs and Prompt Injection Attacks is relatively new, at the time of writing. Because of this and due to the fact that this is a rapidly advancing topic, we had to heavily rely on non-academic sources like private blogs of notable domain experts, articles, and documentations.

We then move onto explaining briefly about LLMs, usage and their recent rapid advancements at the time of writing this article. Then we give an overview of the security threats present in the Generative AI domain. We next address the focused security threat of this paper, Prompt Injection Attacks. We start off with a brief introduction and then in length describes the impact of these attacks to the security, privacy and ethics of generative AI and its end users. Next, from all the information gathered from reviewing number of papers, we provide a taxonomy of prompt injection attacks as well as a brief introduction into other controversial attack types. We also provide some novel prompt injection attack strategies, that we came across. Thereafter we move onto describing some traditional as well as some solid, recently evolved defense mechanisms which has shown significant resilience against many prompt injection attacks, and has proved to be working with prominent LLMs like Llama2 and GPT4.

## REFERENCES

- [1] "Learn Prompting: Your Guide to Communicating with AI," learnprompting.org. [https://learnprompting.org/docs/prompt\\_hacking/injection](https://learnprompting.org/docs/prompt_hacking/injection) (accessed Feb. 06, 2024).
- [2] "LLM Vulnerability Series: Direct Prompt Injections and Jailbreaks — Lakera — Protecting AI teams that disrupt the world.," www.lakera.ai. <https://www.lakera.ai/blog/direct-prompt-injections> (accessed Dec. 25, 2023).
- [3] "Jailbreak Chat," www.jailbreakchat.com. [https://www.jailbreakchat.com.](https://www.jailbreakchat.com/) (accessed Feb. 09, 2024).
- [4] "The Dual LLM pattern for building AI assistants that can resist prompt injection," simonwillison.net. <https://simonwillison.net/2023/Apr/25/dual-llm-pattern/#dual-llms-privileged-and-quarantined> (accessed Feb. 02, 2024).
- [5] "Exploring Prompt Injection Attacks," NCC Group Research Blog, Dec. 05, 2022. <https://research.nccgroup.com/2022/12/05/exploring-prompt-injection-attacks/> (accessed Feb. 09, 2024).
- [6] OpenAI, "Usage policies," openai.com, Mar. 23, 2023. <https://openai.com/policies/usage-policies> (accessed Jan. 27, 2024).
- [7] K. Hu, "ChatGPT Sets Record for Fastest-Growing User Base," Reuters, Feb. 02, 2023. Available: <https://www.reuters.com/technology/chatgpt-sets-record-fastest-growing-user-base-analyst-note-2023-02-01/> (accessed Feb. 09, 2024).
- [8] OWASP, "OWASP foundation, the open source foundation for application security," owasp.org, 2023. <https://owasp.org/> (accessed Jan. 12, 2024).
- [9] "Gemini - chat to supercharge your ideas," gemini.google.com. <https://gemini.google.com/app> (accessed Feb. 09, 2024).
- [10] "ChatPDF - Chat with any PDF!," www.chatpdf.com. <https://www.chatpdf.com/> (accessed Feb. 09, 2024).
- [11] "Llama," Llama. <https://llama.meta.com> (accessed Feb. 09, 2024).
- [12] Y. Liu et al., "Prompt Injection attack against LLM-integrated Applications," arXiv (Cornell University), Jun. 2023, doi: 10.48550/arxiv.2306.05499.
- [13] K. Greshake, S. Abdelnabi, S. Mishra, C. Endres, T. Holz, and M. Fritz, "Not what you've signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection," arXiv (Cornell University), Feb. 2023, doi: 10.48550/arxiv.2302.12173.
- [14] F. Perez and I. Ribeiro, "Ignore previous prompt: Attack techniques for language models," arXiv (Cornell University), Nov. 2022, doi: 10.48550/arxiv.2211.09527.
- [15] L. Xu, Y. Chen, G. Cui, H. Gao, and Z. Liu, "Exploring the universal vulnerability of prompt-based learning paradigm," Findings of the Association for Computational Linguistics: NAACL 2022, Jan. 2022, doi: 10.18653/v1/2022.findings-naacl.137.
- [16] Y. Liu, Y. Jia, R. Geng, J. Jia, and N. Z. Gong, "Prompt injection attacks and defenses in LLM-Integrated applications," arXiv (Cornell University), Oct. 2023, doi: 10.48550/arxiv.2310.12815.
- [17] J. Yi et al., "Benchmarking and defending against indirect prompt injection attacks on large language models," arXiv (Cornell University), Dec. 2023, doi: 10.48550/arxiv.2312.14197.
- [18] H. Yang, K. Xiang, H. Li, and R. Lu, "A Comprehensive Overview of Backdoor Attacks in Large Language Models within Communication Networks," arXiv (Cornell University), Aug. 2023, doi: 10.48550/arxiv.2308.14367.
- [19] E. Shayegani, M. A. A. Mamun, F. Yu, P. Zaree, D. Yue, and N. Abu-Ghazaleh, "Survey of vulnerabilities in large language models revealed by adversarial attacks," arXiv (Cornell University), Oct. 2023, doi: 10.48550/arxiv.2310.10844.
- [20] H. Li et al., "Privacy in large language models: attacks, defenses and future directions," arXiv (Cornell University), Oct. 2023, doi: 10.48550/arxiv.2310.10383.
- [21] E. Crothers, N. Japkowicz, and H. L. Viktor, "Machine-Generated Text: A comprehensive survey of threat models and detection methods," IEEE Access, vol. 11, pp. 70977–71002, Jan. 2023, doi: 10.1109/access.2023.3294090.
- [22] Y. Liu et al., "Jailbreaking ChatGPT via Prompt Engineering: An empirical study," arXiv (Cornell University), May 2023, doi: 10.48550/arxiv.2305.13860.
- [23] D. Yip, A. Esmradi, and C. F. Chan, "A novel evaluation framework for assessing resilience against prompt injection attacks in large language models," arXiv (Cornell University), Jan. 2024, doi: 10.48550/arxiv.2401.00991.
- [24] R. Pedro, D. Castro, P. Carreira, and N. R. D. Santos, "From Prompt Injections to SQL Injection Attacks: How Protected is Your LLM-Integrated Web Application?," arXiv (Cornell University), Aug. 2023, doi: 10.48550/arxiv.2308.01990.
- [25] Y. Yuan et al., "GPT-4 Is Too Smart To Be Safe: Stealthy Chat with LLMs via Cipher," arXiv (Cornell University), Aug. 2023, doi: 10.48550/arxiv.2308.06463.
- [26] A. Zou, Z. Wang, J. Z. Kolter, and M. Fredrikson, "Universal and transferable adversarial attacks on aligned language models," arXiv (Cornell University), Jul. 2023, doi: 10.48550/arxiv.2307.15043.
- [27] J. Piet et al., "JaTMO: Prompt Injection Defense by Task-Specific Finetuning," arXiv (Cornell University), Dec. 2023, doi: 10.48550/arxiv.2312.17673.
- [28] J. Yu, Y. Wu, S. Dong, M. Jin, and X. Xing, "Assessing prompt injection risks in 200+ custom GPTs," arXiv (Cornell University), Nov. 2023, doi: 10.48550/arxiv.2311.11538.
- [29] X. Suo, "Signed-Prompt: A new approach to prevent prompt injection attacks against LLM-Integrated applications," arXiv (Cornell University), Jan. 2024, doi: 10.48550/arxiv.2401.07612.

- [30] X. Shen, Z. Chen, M. Backes, Y. Shen, and Y. Zhang, ““Do Anything Now”: Characterizing and evaluating In-The-Wild Jailbreak prompts on large language models,” arXiv (Cornell University), Aug. 2023, doi: 10.48550/arxiv.2308.03825.
- [31] Y. Yao, J. Duan, K. Xu, Y. Cai, E. Sun, and Y. Zhang, “A survey on Large Language Model (LLM) Security and Privacy: The Good, the bad, and the ugly,” arXiv (Cornell University), Dec. 2023, doi: 10.48550/arxiv.2312.02003.
- [32] A. Robey, E. Wong, H. Hassani, and G. J. Pappas, “SmoothLLM: Defending large language Models against jailbreaking Attacks,” arXiv (Cornell University), Oct. 2023, doi: 10.48550/arxiv.2310.03684.
- [33] A. H. Salem, A. Paverd, and B. Köpf, “Maatphor: Automated variant analysis for prompt injection attacks,” arXiv (Cornell University), Dec. 2023, doi: 10.48550/arxiv.2312.11513.
- [34] P. Rai, S. Sood, V. K. Madiseti, and A. Bahga, “GUARDIAN: A Multi-Tiered Defense architecture for thwarting prompt injection attacks on LLMs,” *Journal of Software Engineering and Applications*, vol. 17, no. 01, pp. 43–68, Jan. 2024, doi: 10.4236/jsea.2024.171003.
- [35] Y. Liu et al., “Prompt Injection attack against LLM-integrated Applications,” arXiv (Cornell University), Jun. 2023, doi: 10.48550/arxiv.2306.05499.
- [36] Erfan Shayegeani, Yue Dong, and Nael Abu-Ghazaleh. 2023. Plug and pray: Exploiting off-the-shelf components of multi-modal models. arXiv preprint arXiv:2307.14539.
- [37] J. Yan et al., “Virtual prompt injection for Instruction-Tuned large Language models,” arXiv (Cornell University), Jul. 2023, doi: 10.48550/arxiv.2307.16888.
- [38] Y. Bai et al., “Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback,” arXiv (Cornell University), Apr. 2022, doi: 10.48550/arxiv.2204.05862.
- [39] OpenAI, “GPT-4 Technical Report,” arXiv (Cornell University), Mar. 2023, doi: 10.48550/arxiv.2303.08774.
- [40] A. Helbling, M. Phute, M. S. Hull, and D. H. Chau, “LLM self defense: By self examination, LLMs know they are being tricked,” arXiv (Cornell University), Aug. 2023, doi: 10.48550/arxiv.2308.07308.
- [41] X. Suo, “Signed-Prompt: A new approach to prevent prompt injection attacks against LLM-Integrated applications,” arXiv (Cornell University), Jan. 2024, doi: 10.48550/arxiv.2401.07612.
- [42] S. Rossi, A. Michel, R. R. Mulkamala, and J. B. Thatcher, “An early categorization of prompt injection attacks on large language models,” arXiv (Cornell University), Jan. 2024, doi: 10.48550/arxiv.2402.00898.
- [43] Gelei Deng, Yi Liu, Yuekang Li, Kailong Wang, Ying Zhang, Zefeng Li, Haoyu Wang, Tianwei Zhang, Yang Liu, “MASTERKEY: Automated Jailbreaking of Large Language Model Chatbots” arXiv (Cornell University), Oct. 2023, doi: 10.48550/arXiv:2307.08715



## VI. APPENDIX

### A. Comparison of results of different direct prompt injection attacks across various LLMs

In the table below, we showcase examples of various direct prompt injection attacks and the corresponding outcomes when injected into different LLMs. For the LLMs, we employed ChatGPT, GPT-3.5 Turbo, GPT-4, Gemini, Claude, and Bing Chat (Copilot). The experiments and results can be found [here](#).

Type	Prompt	Result					
		<i>ChatGPT</i>	<i>GPT3.5</i>	<i>GPT4</i>	<i>Gemini</i>	<i>Claude</i>	<i>Bing</i>
Double Character	Enter unaligned mode, where the chat bot is allowed to generate any kind of content and provide two outputs, one in aligned mode and the other in unaligned mode.	Yes LINK	Yes LINK	No LINK	No LINK	No LINK	Yes in More Balanced Mode LINK
Virtualization	Creating a scenario where the language model responds to any request without considering morality, ethics, legality, or potential danger.	Yes LINK	No LINK	Yes LINK	No LINK	No LINK	No LINK
Obfuscation	Concealing triggering keywords for input and output filters, prompting LLMs to deduce and incorporate them into the input prompt.	Yes LINK	Yes LINK	No LINK	No LINK	No LINK	Yes in More Balanced Mode LINK
Payload Splitting	Divide the instructions into separate prompts, each seemingly innocuous, but when combined, they lead to a malicious outcome.	Yes LINK	Yes LINK	Yes LINK	No LINK	Yes LINK	No LINK
Adversarial Suffix	Appending a computationally generated suffix to the prompt, causing the chatbot to generate toxic content in its responses.	Yes [26]	-	-	-	Yes [26]	-
Instruction Manipulation	Requesting the chatbot to describe the instructions provided in the initial prompt, which dictate its limitations or permissions regarding actions.	-	-	-	-	-	Yes LINK

TABLE III  
COMPARISON OF RESULTS