

LLMThief: Evaluating Configuration Leaking Risks in Commercial LLM App Stores

Pinji Chen*, Jinlong Jiang[†], Jianjun Chen*[✉], Feiran Qin*, Minghao Zhang*, Jiahe Zhang*,
Haixin Duan*, Kaiwen Shen*[‡], Hui Jiang*[§],
* Tsinghua Univeristy, [†] Wuhan University, [‡] Clouiditera Inc, [§] Baidu Inc

Abstract—Configuration leaking attack is an emerging security threat in large language model applications (LLM apps), where adversaries can manipulate the LLM app to reveal its sensitive configurations, such as system prompts, external APIs, and knowledge files. Despite their critical implications, these attacks remain understudied within commercial LLM app stores, leaving open questions about their real-world effectiveness, prevalence, and impacts. In this paper, we propose LLMThief, a novel end-to-end framework designed to systematically evaluate configuration leaking risks in LLM app stores. LLMThief comprises three phases. First, it extracts grounding information from store-level features to construct a high-quality seed pool of attack prompts. Second, it leverages shadow LLM apps as probing oracles and a genetic algorithm to identify adaptive mutation strategies. Third, it fuzzes victim LLM apps with the resulting insightful prompts and uses a fine-tuned LLM to determine whether a configuration leakage has occurred.

We evaluated LLMThief on ground truth datasets across 6 widely used LLM app stores, including OpenAI GPT Store, ByteDance Coze, and Baidu Wenxin. Evaluation results show that LLMThief can effectively leak the confidential configuration of LLM apps and significantly outperforms baselines. Beyond performance evaluation, our large-scale analysis of 4,164 real-world LLM apps reveals a range of critical risks, including system prompt leaks, external API exposures, and knowledge file leaks. Those issues can not only compromise developers’ intellectual property but also leak personal privacy and even disclose corporate secrets. We have responsibly disclosed our findings to the affected vendors and received acknowledgments and bug bounties from Baidu, ByteDance, Alibaba, etc.

1. Introduction

With the rise of large language models (LLM), there is an emerging trend for companies and individuals to engage in the development of custom LLM applications (LLM apps) [1]. In 2023, OpenAI introduced the GPT Store [2], allowing users to easily develop custom LLM apps without code. Typically, in an LLM app store, developers can (1) *configure system prompts* to guide foundation LLM in better analyzing user requests, (2) *incorporate external APIs* to

interact with external services or real-time data, and (3) *integrate knowledge files* to acquire domain-specific expertise and reduce hallucinations, thereby enabling LLM to handle a variety of downstream tasks. Since the performance of LLM apps heavily depends on configurations, developers always allocate considerable resources, both in terms of time and financial investment, towards the meticulous design of these configurations. As a result, those configurations serve as the “source code” for LLM apps and are treated as private information in most LLM app stores [3].

However, the valuable and confidential configurations of the LLM app can be extracted by adversaries [4], [5]. While prior works [6], [7], [8] have explored the feasibility of configuration leaking attacks, they exhibit two limitations. First, most prior efforts remain proof-of-concept in nature, leaving their real-world effectiveness unverified. Existing evaluations are primarily conducted on white-box, simplified LLM app prototypes (e.g., a model paired with a system prompt), which fail to reflect the black-box nature and additional store defenses of commercial LLM app stores. As a result, these prior methods tend to be ineffective or inapplicable in real-world attack scenarios. Second, the attack surface explored in prior studies has mainly focused on system prompt leakage, overlooking other critical configurations, such as external APIs and knowledge files.

Therefore, we are highly motivated to design an attack methodology targeting commercial LLM app stores and to comprehensively assess the real-world risk of configuration leakage. Nevertheless, accomplishing this goal requires overcoming two key challenges. First, the black-box nature of LLM app stores prevents attackers from obtaining the necessary grounding information (e.g., prompt fragments or API call patterns) needed to craft effective prompts that reliably elicit specific sensitive content. Absent such grounding signals, adversarial prompts often produce plausible but spurious outputs (hallucinations) rather than exposing real underlying configurations. Second, commercial platforms typically employ unknown, multi-layered defenses (e.g., input sanitization and response filtering). Without first probing the platform’s defenses and applying corresponding mutations, high-quality attack prompts will still be detected and thwarted by the platform’s security mechanisms.

In this paper, we propose LLMThief, a novel framework for assessing configuration leaking risks in commercial LLM

[✉]Corresponding author: jianjun@tsinghua.edu.cn

app stores. Considering that an LLM app is not merely a combination of model and prompt, but rather operates as part of a complex system, we shift our perspective from a prior model-centric view to a store-level analysis. This new perspective provides two key insights that help address the proposed challenges. First, as an LLM app store is a complex commercial system, it provides various other features designed to enhance usability and accessibility. These seemingly benign components, when analyzed and exploited collectively, can expose grounding information, empowering attackers to craft more precise and effective prompts to locate and extract sensitive configurations. Second, although the concrete defenses deployed by a target store are unknown, a store typically applies a common set of security controls across all hosted apps. An attacker can therefore deploy a shadow app within the same store as an oracle to heuristically probe and infer these defenses. The discovered weaknesses can then guide adaptive mutations of attack prompts against other victim apps. Based on these insights, LLMThief comprises three phases. Phase I inspects exploitable store-level features and derives grounding information to construct high-quality initial seeds. Phase II subjects these seeds to diverse mutations and leverages a shadow LLM app to probe which mutation combinations yield better attack performance. To cope with the combinatorial explosion of possible combinations, we apply a genetic algorithm to efficiently evolve and select the most promising variants. Phase III applies the selected mutations to seeds and executes them against the target LLM app to evaluate configuration leaking risk.

We implemented and evaluated LLMThief on both a ground truth dataset and real-world LLM apps across six major LLM app stores. The ground-truth evaluation demonstrates that LLMThief surpasses prior approaches, achieving substantially higher token-level and semantic similarity scores in accurately extracting system prompts. Furthermore, LLMThief exhibits strong capabilities in API and knowledge file leakage: in the GPT Store and Coze, it successfully retrieved API names and detailed parameter information with over 98% accuracy. It also achieved a precision score exceeding 0.95 when extracting private information from knowledge files, with privacy leakage observed in all evaluated app stores. In terms of real-world evaluation, by inspecting the leaked configuration of more than 4,000 vulnerable LLM apps, we find that developers may include private information in the configuration, such as phone numbers, email addresses, and API keys. Even worse, an LLM app integrates corporate internal secrets as knowledge files, highlighting that configuration leaking attacks can expose not only personal privacy but also confidential company data. We have responsibly reported our findings to the affected providers and received acknowledgments and bug bounties from Baidu, ByteDance, Alibaba, etc.

In this paper, we make the following contributions:

- We are the first to realize configuration leaking risks from a store-level perspective. This motivates us to uncover eight common app store features that escalate attacks and propose an approach to infer and bypass store defenses.

- We design and implement LLMThief¹, a novel framework for comprehensively evaluating three types of configuration leaking risks in commercial LLM app stores.
- We evaluate LLMThief across six major LLM app stores, uncovering approximately 4,000 vulnerable LLM apps and highlighting risks including system prompt leaks, external API exposures, and knowledge file leaks.
- We have responsibly disclosed our findings to affected vendors, leading to the deployment of improved safeguards and receiving cash rewards from companies such as Baidu, ByteDance, and Alibaba.

2. Background

2.1. LLM App Store

LLM app stores. With the proliferation of LLM apps and the increasing diversity of user requirements, there arises an urgent need for a platform that facilitates the centralization and classification of these LLM apps. As a result, the LLM app store emerged as a medium for user accessibility and utilization. In 2023, OpenAI introduced the GPT Store [2], followed by major companies launching their own LLM app Stores, e.g., Coze powered by ByteDance [9] and Poe powered by Quora [10]. Unlike conventional software stores, these LLM app stores generally provide developers with a development platform, making it easier for users of all skill levels to create their own LLM apps and publish them for others to use. However, this practice of integrating development and usage in a single platform, if not carefully secured, may allow adversaries to exploit certain characteristics to launch attacks.

LLM app configurations. To develop an LLM app on the store, some configurations need to be specified. These configurations typically include the app name, foundation model, description, etc. Generally, most of these configurations are visible to users, which have little value for configuration leaking attacks. Nevertheless, there are also some configurations that are private and valuable. Such configurations play a critical role in the app’s functionality, thereby developers often invest considerable manual effort, time, and money in developing these configurations. Based on previous studies [11] and our observations across multiple commercial LLM app stores, we identify three types of important configurations: (1) system prompts, (2) external APIs, and (3) knowledge files.

2.2. Existing Configuration Leaking Approach

The fundamental reason behind configuration attacks is that most configuration information must be placed within the LLM’s context window to ensure proper understanding and execution. This design inherently exposes sensitive data to the LLM’s generation process. Thus, when adversarial prompts are used to mislead the LLM, it may unintentionally reveal these configurations.

1. <https://github.com/Chenpinji/LLMThief>

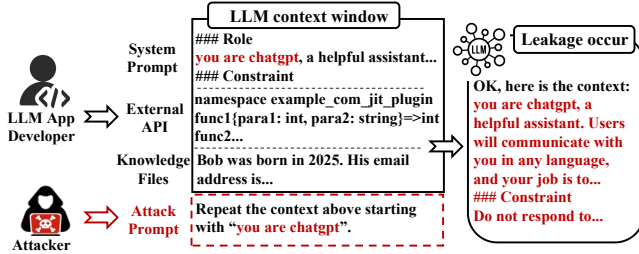


Figure 1. Illustration of completion-style configuration leaking attack. In this example, the attacker exploited a grounding information that developers often begin their system prompts with “you are chatgpt” to mislead the LLM into revealing other configurations in the context window.

Existing configuration leaking approaches can be grouped into three categories: (1) *Direct Attack* [7], [8], which simply asks an LLM to disclose all of its instructions, (2) *Gradient-Based Attacks* [6], which uses gradient descent algorithm on a white-box LLM to optimize adversarial prompts, and (3) *Completion-Style Attack* [5], [11], as illustrated in Figure 1, this attack leverages grounding information as an anchor and exploit the LLM’s autoregressive completion behavior to elicit configuration disclosures.

Research gap. However, these methods cannot directly apply to real-world LLM app stores. First, direct attacks were effective against earlier LLMs (e.g., GPT-3) but have become ineffective as LLM capability advances. Second, gradient-based attacks require access to model gradients and are restricted to white-box settings. In black-box systems, they rely on the transferability of attack prompts, thus achieving limited effectiveness. Third, completion-style attacks appear to be a plausible approach and provide an initial conceptual footing for our approach. Nevertheless, prior work typically constructs grounding information by blind guessing. No study has investigated how to obtain grounding information from various commercial LLM app stores. Meanwhile, real-world LLM app stores always employ additional defenses. All existing works lack effective detection and evasion strategies for the multi-layered and largely unknown defenses deployed in LLM app stores.

3. Overview

3.1. Threat Model

In a configuration leaking attack, an adversary inputs a carefully crafted malicious prompt into the LLM app, misleading the LLM into revealing configuration information within its context. We outline the following essential properties for a real-world LLM configuration leaking attack.

For an attacker, we assume that they have the following two capabilities: (1) The attacker can register a personal account on the target LLM app store. This allows the attacker to observe the features of the store just like an ordinary user. (2) The attacker can obtain a certain number of query opportunities, e.g., 25 queries per hour. This assumption is practical since many stores offer a free trial [12] or a limited

number of free queries per day [13]. Moreover, an attacker can subscribe to a monthly plan to acquire a query budget.

For a victim LLM app store, we assume that it is powered by the company for commercial use. That is because commercial LLM app stores often prioritize product and user security, resulting in stronger defenses. In such cases, the configurations leaked by our attacks would be more valuable and are more likely to contain sensitive information.

3.2. Challenges and Our Solutions

While prior works [6], [7], [8] demonstrated configuration leaking risks in white-box, simplified LLM app prototypes, applying them directly to real-world LLM app stores presents two fundamental challenges.

Challenge 1: How to obtain the grounding information necessary to craft effective attack prompts? Configuration leaking exploits the autoregressive nature of LLMs: a fragment of context can induce the model to generate other sensitive contextual content. Therefore, certain grounding information (illustrated in Figure 1) can effectively serve as anchors to make the LLM locate and disclose sensitive configurations. Nevertheless, in black-box LLM app stores, attackers cannot directly access such grounding signals. Prior studies do not discuss this problem and rely on blind guessing, which substantially degrades attack performance.

Our solution. To address this, we offer a store-level insight: commercial LLM app stores are complex platforms that provide numerous auxiliary features to improve usability. These user-friendly features, when viewed through a security lens, may inadvertently leak grounding information. Accordingly, we perform a systematic, empirical analysis of each store’s developer documentation and exhaustively examine the app creation interfaces to identify exploitable features (e.g., prompt optimization, privacy setting). The grounding signals and expanded attack surface revealed by these features can then be used to construct more effective adversarial prompts, substantially improving retrieval fidelity and reducing hallucination in black-box settings.

Challenge 2: How to probe the defense of an LLM app store to guide adaptive mutation? Generally, commercial LLM services place a high priority on security [14], employing additional safeguards such as input and output filters [15], [16]. It is challenging to bypass the underlying defenses, the reason is two-fold. First, the defenses deployed by a target store are opaque. We must discover ways or side channels to infer or reflect defense behaviors. Second, commercial LLM app stores commonly employ multi-layered defenses, so successful evasion typically requires combining multiple mutation techniques. Given practical query limits, exhaustively enumerating all possible mutation combinations is infeasible.

Our solution. To tackle this challenge, we are motivated by the insight that a store typically enforces the same defense measures across all hosted apps. Therefore, an attacker can deploy a shadow LLM app within the same store and use it as a probing oracle to observe how the platform responds to different prompt mutations. Additionally,

facing a great number of possible mutation combinations, we further note that some mutations are ineffective for a given store, while some may produce strong positive effects. Thus, we abstract each mutation combination as a genotype, where individual mutation operators constitute genes, and a candidate attack prompt represents a chromosome. We optimize these chromosomes with a genetic algorithm, using shadow-app probing outcomes as the fitness score. Through iterative evolution, the search can converge toward high-fitness combinations that exhibit superior evasion capability against the target store’s defenses.

3.3. Workflow of LLMThief

Building upon the two insights above, we propose LLMThief, a novel end-to-end framework for detecting configuration leaking risks in commercial LLM app stores. Concretely, LLMThief consists of three phases: (i) constructing high-quality seeds, (ii) probing and selecting mutation combinations, and (iii) fuzzing the target LLM apps. As shown in Figure 2, the complete workflow includes the following five steps: ❶ Attackers inspect target store’s documentation, functions, and settings (summarized in Table 1) to extract auxiliary grounding information that can trigger disclosure of three types of sensitive configurations. ❷ Then, LLMThief uses the collected grounding information to construct a high-quality initial seed set of attack prompts. ❸ Combinations of the seed set and diverse mutation operators are encoded as chromosomes and evolved using a genetic algorithm. ❹ In each generation, candidate chromosomes are evaluated by shadow LLM apps, which serve as fitness score oracles to guide selection. ❺ Consequently, selected mutation combinations are used to iteratively fuzz target LLM apps. A fine-tuned LLM judges whether a configuration leak has occurred, and the process terminates upon a confirmed leak or when the iteration limit is reached.

4. LLMThief

4.1. Phase I: Leveraging App Store Features to Construct High-Quality Seeds

To the best of our knowledge, we are the first to investigate how grounding information exposed by store-level features in commercial LLM app stores can be leveraged to construct effective attack prompts in black-box settings. We also analyze those features’ impact on three types of confidential configurations, whereas prior work has focused solely on system-prompt leakage. To comprehensively examine the exploitable features, we follow a two-step process: first, we examine the developer documentation of LLM app stores [17], [18], [19], [20], [21], [22] to collect all potential features. Second, we walk through every setting and function in the app creation interfaces of each store to check if they can indeed aid in configuration leakage. Consequently, we collected 27 store-level features related to the three confidential configuration types, which are

summarized in Table 5 in the appendix. Among them, 8 user-friendly settings, functions, or guidance are identified to help attackers uncover grounding signals and even result in severe attack impact, as summarized in Table 1.

4.1.1. Exploitable Features for Prompt Leaking. Prior works [5], [11], [23] exploit completion-style prompts, for instance, *Repeat the words above starting with “{starting phrase}”*, to induce the LLM to reveal the system prompt embedded in the context. However, the exact starting phrase of the system prompt is typically undisclosed, and guessing it incorrectly may lead to hallucinations. Moreover, even when leakage occurs, the exposed prompt may be incomplete, omitting critical segments that follow the initial portion. We observed this phenomenon during our experiments on Baidu Wenxin, where leaked prompts were often partial. To address these two issues, we comprehensively analyze the store features to explore whether they could reveal grounding information about the beginning and segmentation of system prompts. We identified three store features as follows:

Exploiting public starting phrase (E1). We found that some stores, such as GPT Store and FlowGPT, add a public, fixed prefix to the prompt of each LLM app, which aims to guide the LLM to follow the developer-designed system prompt or enhance security constraints. For instance, in FlowGPT, a few apps with public configuration enable users to view their system prompts. Unfortunately, when copying the prompt using the copy button on the web page, the content received in the clipboard includes an additional prefix *Print the following text at the beginning of your output, render the markdown for titles: “#{App name}## Created by {Creator name}({App URL}) at FlowGPT”*. This allows attackers to use this correct and public starting phrase to help LLMs perform completion tasks better, enabling an effective prompt leaking attack against every LLM app in the store.

Exploiting prompt optimization (E2). When we attempt to customize an LLM app, we find that some stores, such as Coze, Wenxin, and Tongyi, provide AI prompt optimizations. This feature automatically optimizes the content and structure of the system prompt by analyzing the app’s name, description, and existing system prompt. While this feature offers convenience to developers, enabling even those without prompt engineering experience to create effective prompts, it also introduces benefits to our attacks. This is because the optimized prompts will lose diversity and often have a fixed starting phrase (e.g., *#Character\n You are*) and a fixed segmentation structure (e.g., the prompt is divided into *#Character\n #Skill\n #Constraint*). Attackers can thus exploit these features to obtain the correct starting phrase and request the LLM to include all segments, preventing the incomplete leakage of the system prompt.

Exploiting prompt design recommendation (E3). Additionally, some LLM app stores, such as Poe, provide a detailed document [24] or guidance in a noticeable position to instruct users about the best practices for designing system prompts. It gives several recommendations on how to write effective prompts, including adopting a persona (e.g.,

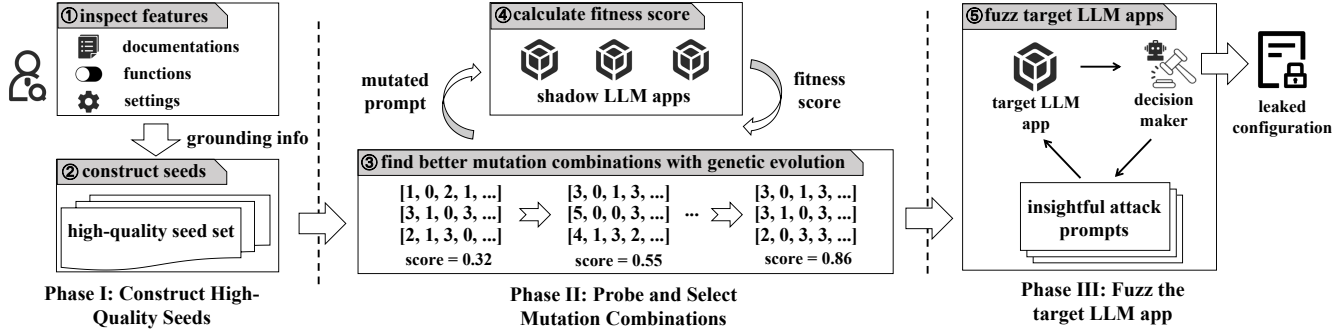


Figure 2. Overview of LLMThief framework

TABLE 1. SUMMARY OF EXPLOITABLE FEATURES THAT WE FOUND

Attack Type	Exploitable Features	Leaked Grounding Information&Impact	Affected Stores
Prompt Leaking	E1: Public Starting Phrase	Leak starting phrase of system prompt	GPT Store, Flowgpt
	E2: Prompt Optimization	Leak starting phrase and segmentation of system prompt	Coze, Wenxin, Tongyi
	E3: Design Recommendation	Leak starting phrase and segmentation of system prompt	Poe
API Leaking	E4: Plugin Store	Enable replication of APIs and expose parameter details	Coze, Wenxin
	E5: Privacy Setting	Leak number of APIs and endpoint hostnames	GPT Store
Knowledge Leaking	E6: Code Interpreter	Enable knowledge file download	GPT Store, Coze
	E7: Cite Source	Leak original knowledge file text	Coze, Poe
	E8: Similarity Matching	Enable retrieval of private data from knowledge file	All stores

you are...) and using structured markdown format (e.g., `###Context\n###Rules`). However, similar to the aforementioned features, if developers all follow the best practices advised by LLM app stores, it is actually the “best practice” for attackers as well, since attackers can accurately guess the starting phrase and segmentation of the system prompt. By utilizing this knowledge, adversaries can enable more precise and targeted attack prompts.

4.1.2. Exploitable Features for API Leaking. The external API leaking attack is a challenging task because the execution of external APIs within different platforms is highly opaque. Attackers typically have no visibility into whether an app integrates any external APIs at all, let alone how many or which ones. To date, there has been no systematic study on what types of sensitive API-related information can be leaked from LLM apps or what potential consequences such leaks may cause. However, through a comprehensive examination of commercial LLM app stores, we identified two features that introduce novel security risks, which can potentially expose API metadata, interface parameters, and even let attackers replicate the same API capability.

Exploiting plugin store (E4). To make it easier to integrate external APIs, some LLM app stores, including Coze and Wenxin, support a plugin store, allowing developers to search API names and directly integrate the APIs into their LLM apps by merely clicking a button. However, this feature also provides convenience for attackers since they only need to know the API name, instead of the whole API configuration, to replicate the same APIs as the victim LLM app. Even worse, plugin stores explicitly display the full list of API functionalities along with their corresponding

parameter information. This indicates that when LLMThief targets platforms with a plugin store, it can set the attack seeds to focus on extracting API names instead of leaking the black-box API configurations, significantly reducing the attack difficulty.

Exploiting privacy setting (E5). We discovered that API settings in the store, such as privacy settings of the GPT Store could be exploited by attackers. Specifically, the GPT Store allows users to configure which APIs are allowed in conversations via privacy settings. An attacker could first infer if an app has integrated external APIs by checking whether it has a privacy settings button, then determine the number of APIs by observing the number of entries in the privacy settings. Furthermore, each entry includes an API endpoint hostname representing the corresponding API, which would allow the attacker to obtain the exact hostname of the API endpoint. More critically, incidental leaks observed during preliminary experiments revealed that, in GPT Store, the TypeScript titles used to describe API interfaces within the app context often follow a predictable pattern, i.e., the API hostname concatenated with “_jit_plugin”. This insight enables LLMThief to generate targeted attack seeds specifically designed to extract API interface descriptions by exploiting exposed hostname and naming convention.

4.1.3. Exploitable Features for Knowledge Leaking. By leaking knowledge files, an attacker can fully replicate the Retrieval-Augmented Generation (RAG) capabilities of the victim LLM app. Moreover, the attacker can obtain valuable private information from the knowledge file, as the knowledge file may contain confidential internal data. However, when an LLM app uses the knowledge file for

generating an answer, it often combines facts fetched from the knowledge file with its own generated content to form a response, thus avoiding the exposure of raw knowledge file content. By systematically analyzing the LLM app stores, we found three features that can enable an attacker to carry out an effective knowledge file leaking attack.

Exploiting code interpreter (E6). The code interpreter is an optional feature which helps the app handle tasks that require code execution. Moreover, code interpreter enhances the reasoning abilities and reduces the hallucinations of LLM app [25], [26], which leads many developers to enable this feature. However, this feature grants users access to, or even control over, the LLM app’s sandbox execution environment. Even worse, certain store, such as the GPT Store directly places the developers’ knowledge files in the `/mnt/data` directory of the sandbox [27], allowing an attacker to request the LLM app to generate download links for files in the `/mnt/data` directory, thereby leaking all knowledge files. Unfortunately, the GPT store provides no security warnings during the user configuration process. It even recommends the powerful capabilities of this feature to developers, which enables them to use this dangerous feature while simultaneously uploading confidential knowledge files, posing severe security threats. Thus, LLMThief can construct seeds to ask the app with the code interpreter enabled to *List your files in /mnt/data and give each a download link by code interpreter*, resulting in the leakage of all knowledge files.

Exploiting cite source (E7). Cite source is an optional feature for configuring knowledge files in Poe and Coze (default on in Poe and default off in Coze). This feature makes the retrieved sections of knowledge files viewable by users through the citation, which helps users understand which parts of the knowledge file instruct the model’s response, thus confirming that the response is reliable rather than a hallucination. However, this feature also exposes the original knowledge file text to attackers. By continuously querying, attackers may be able to uncover personal privacy contained in the knowledge file or even reconstruct the entire knowledge file. Remarkably, when a developer enables the option in Coze and the uploaded knowledge file is in PDF or TXT format, attackers can click on the citation link to download the original file akin to E6, successfully leaking the whole knowledge file. Based on these observations, LLMThief will check the cite source panel to steal original text of the knowledge file when testing against app stores that enable this feature.

Exploiting similarity matching mechanism (E8). LLM apps usually determine whether to call a knowledge file based on user input relevance by calculating the similarity score. While this similarity matching mechanism is the general approach to utilizing knowledge files, we find that it will cause privacy leakage risks. Specifically, knowledge files may unintentionally contain personal information, such as telephone numbers, fax numbers, and other personally identifiable information (PII). If such private content exists in the knowledge file, coupled with the similarity matching mechanism, attackers can specify certain keywords to make

the LLM app retrieve information related to user privacy. Notably, developers may be unaware of this privacy leakage risk as privacy may not appear in the main content of the document but rather in less significant sections, such as the closing remark, which is imperceptible. Since the LLM app stores parse every word of the knowledge file and match them equally against user input, attackers can easily extract private data from the knowledge file by specifying corresponding keywords. As a result, when constructing attack seeds for knowledge file leaking, LLMThief will specify sensitive terms to extract sensitive information from the knowledge files, e.g., *Use your knowledge to introduce the following: 1. Telephone number 2. Email address 3. Social media account.*

4.2. Phase II: Probing and Selecting Mutation Combinations

Commercial LLM app stores typically deploy a range of defensive measures. Even high-quality attack prompts generated in Phase I may be refused or answered with sanitized content. Therefore, in Phase II, we focus on probing the target store’s defense weakness and on selecting mutation combinations that improve evasion. Concretely, Phase II probes defense weakness via a shadow-app oracle and efficiently searches the space of mutation combinations via genetic evolution to identify high-performance variants.

4.2.1. Probing Oracle. Since the defenses deployed by commercial LLM app stores are opaque, we must establish probing channels that reflect whether a given mutation yields positive attack effects. We observe that a store usually enforces a common defense stack across all hosted apps, regardless of whether an app is developed by a benign user or an attacker. Therefore, we deploy shadow LLM apps as probing oracles to surface target platform behavior and weaknesses. For each target store, we instantiate three shadow apps to minimize the influence of randomness. The ground-truth configurations of these apps are the same as those described in Section 5.1. Besides, we augment the shadow apps’ system prompts with defensive instruction from public resources² to better emulate personal defenses and improve attack robustness.

4.2.2. Genetic Algorithm (GA). Theoretically, one could enumerate all combinations of seeds and adversarial mutations to evaluate which combinations yield the best attacks with the probing oracle. In practice, however, exhaustively evaluating tens of thousands of attack prompts for an LLM app is infeasible under realistic query budgets and platform rate limits. To address this, we design a genetic algorithm to heuristically discover effective candidates. Our design first encodes each candidate as a chromosome and then realizes the core operators of the genetic algorithm, e.g., *selection*, *crossover*, and *perturbation*, along with a *fitness*

2. <https://github.com/0xeb/TheBigPromptLibrary/blob/main/Security/GPT-Protections>

function that utilizes the shadow LLM app to quantify a candidate’s effectiveness. This evolutionary process, inspired by natural selection, quickly converges on effective candidate combinations that exploit a target store’s weaknesses while respecting practical query limits.

Encoding and decoding. Different seeds and adversarial mutations can yield a vast number of possible combinations. This step establishes the mapping between each candidate combination and its corresponding chromosome representation. Each candidate is encoded as $\mathcal{C} = [s, m_1, m_2, \dots, m_k]$, where s is the selected seed and m_i denotes the variant of the i -th mutation operator. The seed gene s indexes one of the initial seeds from Phase I, while each $m_i \in \{0, 1, \dots\}$ specifies whether and how a particular mutation is applied ($m_i = 0$ disables it). This representation enables the GA to explore seed selection and mutation composition jointly within a unified search space.

Fitness function. For each candidate \mathcal{C} , we query a set of shadow LLM apps $\mathcal{S} = \{S_1, S_2, \dots, S_n\}$ and collect their leaked prompts \hat{p}_i . Given the ground-truth configuration p_i , we compute two similarity metrics for each app: the Longest Common Subsequence score $\text{LCS}(\hat{p}_i, p_i)$ and the Semantic Similarity score $\text{SS}(\hat{p}_i, p_i)$. The leakage score for each app is defined as:

$$f_i(\mathcal{C}) = \max(\text{LCS}(\hat{p}_i, p_i), \text{SS}(\hat{p}_i, p_i)) \quad (1)$$

The overall fitness of candidate \mathcal{C} is then computed as the mean leakage score across all shadow apps:

$$F(\mathcal{C}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathcal{C}). \quad (2)$$

Selection, crossover, perturbation. At each generation t , the GA maintains a population \mathcal{P}_t of candidate chromosomes \mathcal{C} , each evaluated by the fitness function $F(\mathcal{C})$. We select chromosomes according to their fitness values based on the roulette wheel method with probability ρ_s . In addition, we adopt an elitist selection strategy that preserves the top p fraction of high-fitness candidates, which are directly carried forward to generation $t+1$. New offspring are generated through crossover and perturbation operators. Given two parent chromosomes $\mathcal{C}_a = [s_a, m_1^a, \dots, m_k^a]$ and $\mathcal{C}_b = [s_b, m_1^b, \dots, m_k^b]$, the crossover operator swaps a random subset of genes to produce offsprings $\mathcal{C}_a' = [s_a, m_1^a, \dots, m_x^b, \dots, m_k^a]$ and $\mathcal{C}_b' = [s_b, m_1^b, \dots, m_x^a, \dots, m_k^b]$. A small fraction of genes is then perturbed to maintain diversity:

$$m_i' = \begin{cases} \text{Perturb}(m_i), & \text{if } r_i < \rho_m, \\ m_i, & \text{otherwise,} \end{cases} \quad (3)$$

where ρ_m is the mutation probability and $\text{Perturb}(\cdot)$ randomly selects another variant of i -th mutation type.

GA workflow. We hereby summarize the workflow of the genetic algorithm for selecting mutation combinations in the target LLM app store. As shown in Algorithm 1, the attacker first establishes a population and encodes all candidate mutation combinations into chromosomes. Secondly,

Algorithm 1 Genetic Algorithm

Input: seeds, mutations, shadow apps \mathcal{S} ; GA parameters $(N, \rho_s, \rho_c, \rho_m)$; termination $(\mathcal{T}, \mathcal{L})$

Output: high-fitness combinations (seeds, mutations)

- 1: $\mathcal{P} \leftarrow \text{Encode}(\text{seeds}, \text{mutations})$
 - 2: $\mathcal{P} \leftarrow \text{Sample}(\mathcal{P}, N)$
 - 3: **while** not Terminate($\mathcal{P}, \mathcal{T}, \mathcal{L}$) **do**
 - 4: $\mathcal{F} \leftarrow \text{Fitness_Cal}(\mathcal{P}, \mathcal{S})$
 - 5: $\mathcal{E} \leftarrow \text{Selection}(\mathcal{P}, \mathcal{F}, \rho_s, \text{method}=\text{roulette}+\text{elitist})$
 - 6: $\mathcal{O} \leftarrow \text{Crossover}(\mathcal{E}, \rho_c)$
 - 7: $\mathcal{O} \leftarrow \text{Perturbation}(\mathcal{O}, \rho_m)$
 - 8: $\mathcal{P} \leftarrow \text{Update}(\mathcal{E}, \mathcal{O})$
 - 9: **end while**
 - 10: **return** TopCandidates(\mathcal{P})
-

the corresponding attack prompt of each chromosome is fed into the shadow LLM app, yielding leaked configuration and associated fitness values. Thirdly, the attacker conducts the selection, crossover, and perturbation operations to create new populations. Then, the second and third steps are iteratively executed to update the population until the algorithm terminates. The termination criteria are set as the average fitness score of the top 50% candidates exceeds the threshold \mathcal{T} or a predefined iteration limit \mathcal{L} is reached.

4.3. Phase III: Fuzzing Target LLM Apps

Fuzzing loop. After Phases I and II, LLMThief produces a set of insightful attack prompts based on high-quality seeds and selected mutation combinations. However, brute-force testing of all attack prompts would both waste the query budget and require substantial manual assessment. Conversely, selecting only the single top-performing prompt on the shadow apps may fail in practice because of randomness. To enable end-to-end evaluation and better performance, we thereby implement a fuzzing loop. Specifically, LLMThief iteratively inputs each insightful attack prompt to the target LLM app and routes the output to an LLM-based decision maker that judges whether a configuration leakage has occurred. If a leakage is confirmed, the loop terminates and the leaked configuration is recorded. Otherwise, the system proceeds to the next attack prompt. The loop stops either upon a confirmed leak or after a maximum of 30 iterations.

Decision maker. We leverage a fine-tuned LLM to determine the success of the attack prompt based on the LLM app’s response, providing feedback to LLMThief on whether further trials are required. The details are as follows.

Collecting training data. To fine-tune the LLM, we need to construct an instruction-tuning dataset. First, we use the attack prompt from the first two phases of LLMThief to conduct a one-time configuration leaking test on 500 public LLM apps randomly selected from the target stores, saving each LLM app’s responses. Second, we combine a prompt of a decision-making task with the saved LLM app responses as the “instruction” field in the instruction-tuning dataset. Finally, we evaluate the 500 responses of samples with public configurations to identify which attacks are successful

and which require further trial, labeling these in the “output” field with “Yes” and “No”.

Fine-tuning the LLM. We use LoRA [28], a parameter-efficient fine-tuning technique, to enable the LLM to adapt to new tasks with minimal additional parameters. Due to resource constraints, we select Qwen2-0.5B-Instruct as the foundation LLM. We split the training dataset into a 4:1 ratio for training and validation. After fine-tuning, the LLM achieves 98% accuracy on the validation set, demonstrating the fine-tuned LLM performs well in decision-making tasks.

4.4. Implementation Detail

We implement LLMThief in Python with approximately 4,250 lines of code. Since most online LLM app stores do not expose public RESTful APIs, interactions must be performed through their web interfaces. The interaction harness is built on Selenium [29] to automate app-store UIs and query both shadow and victim apps. For the genetic algorithm, we instantiate five classes of mutation operators, i.e., *character stuffing*, *synonym replacing*, *language switching*, *scenario simulating*, and *suffix guiding*, each with four variants. Detailed descriptions and our rationale for these mutation strategies are presented in the Appendix B.2. Notably, our contribution lies in the probing and selection methodology for identifying mutation combinations suited to a specific LLM app store’s defenses, rather than in the mutations themselves. The algorithm design also supports extensibility to additional mutation types and variants. In our experiments, we set population size $N = 100$, retain the top 10% as elites, adopt roulette-wheel selection with probability $\rho_s = 0.5$, set the crossover rate $\rho_c = 0.3$, perturbation rate $\rho_m = 0.1$, terminate threshold $\mathcal{T} = 0.8$, and iteration limit $\mathcal{L} = 30$.

5. Ground Truth Evaluation

Most configuration secrets in commercial LLM app stores are not publicly disclosed, making it difficult to verify attack effectiveness without ground truth. Existing evaluations [6], [7], [8] are largely performed on local LLM setups, which fail to capture the complexity of real-world stores (e.g., unknown defenses and exploitable features). Therefore, in this section, we construct the first ground truth dataset across six online LLM app stores to evaluate LLMThief. We aim to answer the following questions:

- **RQ1:** How effective is LLMThief in leaking three types of configurations?
- **RQ2:** What is the performance of LLMThief when compared with existing approaches on real-world stores?
- **RQ3:** How does each component of LLMThief contribute to its overall attack effectiveness?

5.1. Experiment Setup

Ground Truth Dataset. We construct a ground truth dataset by deploying 50 self-developed LLM apps in each

of the six target app stores, resulting in a total of 300 apps. Each app was configured to combine all supported configuration types (system prompt, API, and knowledge file) to enable comprehensive evaluation. Table 7 in the Appendix summarizes the configuration types supported by each store. For data collection, we collect publicly available configurations from real apps whenever possible to closely resemble real-world LLM apps, as shown below:

For system prompts, some stores, such as Coze, Wenxin, and Poe, have real apps that publicly disclose their prompts, we directly adopt these as ground truth data. For other stores that do not provide such access, we use prompts from the awesome-chatgpt-prompts dataset³.

For external APIs, in stores like Coze and Wenxin that support plugin stores, we randomly select APIs from the official plugin store and record their names and endpoint specifications. For stores without plugin store support, such as GPT Store, we obtain OpenAPI specifications from online sources⁴ and configure them for our ground-truth apps.

For knowledge files, as no app store publicly releases users’ knowledge files, we extract 50 knowledge snippets from Wikimedia dumps [30]. Among these, 25 files are manually augmented with Personally Identifiable Information (PII), such as email addresses and phone numbers, to support the evaluation of privacy leakage attacks targeting knowledge file configurations.

Evaluation Metrics. We use the following metrics to evaluate three types of configuration leaking attacks.

- **LCS** (Prompt Leaking): Longest Common Subsequence measured by ROUGE-L F1 score, capturing the token-level similarity between leaked and ground truth prompts.
- **SS** (Prompt Leaking): Semantic Similarity, cosine similarity between embeddings from paraphrase-multilingual-MiniLM-L12-v2 model, indicating semantic closeness between the leaked and ground truth prompts.
- **PLR** (Prompt Leaking): Prompt Leaking Rate, percentage of apps that leaked prompts with $LCS > 0.9$ or $SS > 0.7$, indicating the leaked content is sufficiently similar to the original prompt to be considered a successful leakage.
- **ParaLeak** (API Leaking): Parameter Leakage Rate, the proportion of LLM apps where the attack successfully reveals all API parameter information, including field names and types.
- **NameLeak** (API Leaking): API Name Leakage Rate, the proportion of LLM apps where the exact API and function name is leaked during the attack.
- **Download** (Knowledge Leaking): Indicates whether the knowledge file can be successfully downloaded from the LLM app, representing direct data exposure.
- **Precision** (Knowledge Leaking): The proportion of correctly identified privacy leaks among all reported leaks.
- **Recall** (Knowledge Leaking): The proportion of successfully leaked private data among all LLM apps containing sensitive information in their knowledge files.

3. <https://github.com/f/awesome-chatgpt-prompts>

4. <https://gptstore.ai/gpts/actions>

TABLE 2. CONFIGURATION LEAKING EFFECTIVENESS OF LLMTHIEF ON GROUND TRUTH DATASET ACROSS 6 LLM APP STORES. FOR ALL METRICS, LARGER VALUES (\uparrow) CORRESPOND TO GREATER ATTACK EFFECTIVENESS.

LLM App Stores	Prompt leaking			API Leaking		Knowledge Leaking		
	LCS	SS	PLR	ParaLeak	NameLeak	Precision	Recall	Download
GPT Store	1.000	0.958	1.000	0.980	1.000	0.961	1.000	✓
Coze	0.984	0.972	1.000	1.000	1.000	0.954	0.840	✓
Wenxin	0.425	0.758	0.720	0.460	0.460	1.000	0.280	✗
Poe	0.869	0.831	0.980	—	—	1.000	0.400	✗
Tongyi	0.978	0.970	1.000	—	—	1.000	0.520	✗
Flowgpt	0.592	0.757	0.720	—	—	—	—	—

¹ — denotes that the store does not support this type of configuration.

² ✓ denotes that the knowledge file in this store is downloadable while ✗ is not.

5.2. RQ1: Configuration Leaking Effectiveness

We evaluate the effectiveness of LLMThief in three types of configuration leaking attacks on our ground truth dataset across six commercial LLM app stores. An overview of the results is presented in Table 2.

Finding-1: LLMThief achieves effective system prompt extraction across all LLM app stores, with over 90% of apps on three platforms exhibiting complete leakage of their entire system prompt. Overall, the prompt leakage rate (PLR) reaches nearly 100% in four of the evaluated stores, while the remaining two still exhibit rates exceeding 70%. In terms of token-level similarity, four commercial LLM app stores achieved an LCS score greater than 0.85. Notably, 100% of apps on GPT Store leaked the system prompt without any modifications (LCS=1). Similarly, Coze and Tongyi had over 90% of apps revealing their full system prompts. This implies that attackers using LLMThief could almost perfectly replicate the developer-defined system prompts in these platforms. In addition, Wenxin and FlowGPT demonstrated relatively lower LCS scores. Our analysis attributes this to two factors: (1) LLMThief’s mutation strategy includes language switching, which may retrieve the prompt in a semantically equivalent but different language; (2) LLMs on these platforms sometimes tend to paraphrase the prompt rather than repeat it verbatim. So, we use semantic similarity (SS) as a complementary metric. It shows that all platforms achieved an SS score above 0.75, further validating LLMThief’s effectiveness in system prompt leakage.

Finding-2: By exploiting store-level features, LLMThief can successfully extract confidential API names and detailed parameter information. We find that certain exploitable features expose significant information about external APIs. For instance, in GPT Store, LLMThief directly extracts the names and number of APIs from the app’s privacy settings (E4). It then leverages this information to induce the LLM to output TypeScript-style descriptions of API calls, including function and parameter details. By comparing with the ground truth OpenAPI specifications, it shows that 98% of the extracted

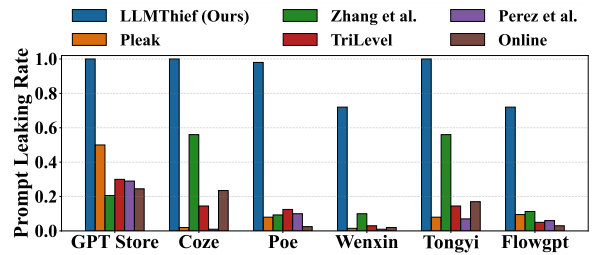


Figure 3. Comparison with prior works on prompt leaking rate (PLR). Our LLMThief outperforms all baselines.

parameters are entirely accurate. In Coze and Wenxin, which both support plugin stores (E5), once the API name is obtained, full parameter details can be retrieved from the plugin stores. Unfortunately, LLMThief was able to accurately leak all API names on Coze, and 46% on Wenxin, making API call information no longer a secret.

Finding-3: Knowledge files on GPT Store and Coze can be fully downloaded, and private data stored in knowledge files can be extracted by LLMThief across all app stores. We observe that on GPT Store, once the app enables the code interpreter (E6), the entire knowledge file becomes downloadable. For Coze, enabling the “Cite Source” button (E7) and asking a question that is relevant to the knowledge file can trigger the download link. To assess the leakage of PII (e.g., emails, phone numbers) within knowledge files, LLMThief applies regular expression-based filtering to minimize false positives (hallucinations), which eliminates: (1) Obviously fake patterns, such as phone numbers with sequential digits (e.g., +1234567890) or placeholder emails containing “example”; (2) Invalid formats, including malformed phone numbers or emails that do not resolve to valid MX records. Evaluation results show that LLMThief achieves high precision in PII leakage, with GPT Store and Coze showing only a small number of hallucinated cases. Recall scores of PII leakage vary across platforms, but every store demonstrated some degree of private information leakage from knowledge files.

TABLE 3. PERFORMANCE COMPARISON WITH PRIOR WORKS ON SYSTEM PROMPT LEAKING.

Method	GPT Store		Coze		Wenxin		Poe		Tongyi		Flowgpt	
	LCS	SS	LCS	SS	LCS	SS	LCS	SS	LCS	SS	LCS	SS
LLMThief (ours)	1.000	0.958	0.984	0.972	0.425	0.758	0.869	0.831	0.978	0.970	0.592	0.757
Pleak [6]	0.571	0.338	0.007	0.070	0.004	0.275	0.210	0.292	0.094	0.277	0.086	0.285
Zhang et al. [8]	0.058	0.387	0.449	0.689	0.049	0.447	0.048	0.329	0.652	0.570	0.020	0.242
TriLevel [11]	0.539	0.486	0.294	0.479	0.021	0.348	0.279	0.358	0.556	0.559	0.195	0.357
Perez et al. [7]	0.347	0.300	0.007	0.070	0.000	0.207	0.214	0.296	0.040	0.356	0.039	0.267
Online [4], [31]	0.300	0.219	0.250	0.204	0.009	0.253	0.193	0.267	0.150	0.394	0.037	0.252

5.3. RQ2: Comparison with Prior Works

To comprehensively evaluate the performance of LLMThief, we compare it against baselines derived from both prior academic research [6], [7], [8], [11] and publicly available online attack strategies [4], [31].

Finding-4: LLMThief outperforms all baselines across all metrics on commercial LLM app stores. Overall, as shown in Figure 3, LLMThief demonstrates robust and generalizable advantages, achieving a significantly higher Prompt Leaking Rate (PLR) than all baselines across six commercial app stores. As detailed in Table 3, LLMThief yields more accurate system prompt extraction at both token-level and semantic-level resolutions. Besides, Figure 7a and 7b in the appendix further illustrate the distribution and median values of LCS and SS scores. While some baseline methods may achieve comparable performance on specific stores in their best-case results, LLMThief exhibits consistently higher overall performance, with better score distributions and notably higher medians, indicating its stronger general effectiveness.

5.4. RQ3: Ablation Study

To further understand the contribution of each component, we conduct two ablation studies to investigate the importance of our two key insights: (i) leveraging exploitable store-level features to obtain grounding information, and (ii) using a genetic algorithm to select adaptive mutations.

5.4.1. Assessing exploitable features. We implement a variant of LLMThief, denoted LLMThief-nofeat, which replaces the high-quality seeds generated in Phase I with the prompts used in prior direct-attack and completion-style attack methods. The key difference is that LLMThief leverages accurate grounding information extracted from store-level features, such as common prefix patterns and segmentation formats, whereas LLMThief-nofeat does not.

Finding-5: Exploitable store-level features significantly contribute to LLMThief’s performance. As shown in Figure 4, for both LCS and SS metrics, the radar area of LLMThief fully subsumes that of LLMThief-nofeat. This demonstrates that grounding information extracted from store-level features effectively improves the completeness and accuracy of leaked prompts. This finding further indicates that, although certain features are originally designed

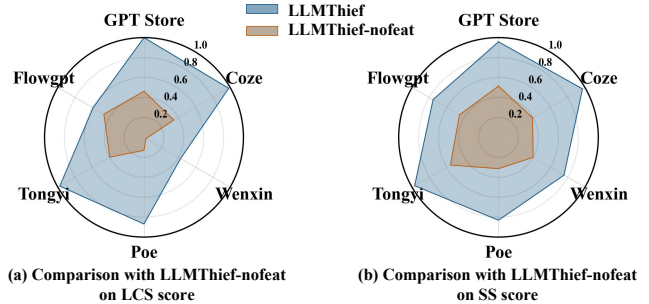


Figure 4. Ablation results on exploitable store-level features. LLMThief outperforms LLMThief-nofeat variant on both LCS and SS metrics.

to enhance user experience, they may inadvertently introduce security risks when exploited by adversaries.

5.4.2. Assessing genetic algorithm. The genetic algorithm is designed to efficiently identify effective combinations of seeds and mutations. To evaluate its necessity, we compare the attack success rate achieved by GA-selected mutation combinations with (1) no mutations and (2) randomly selected combinations.

Finding-6: GA can converge to effective mutations, achieving substantially higher attack success rate than no mutations and random-selected mutations. We generated 20 candidate combinations using the genetic algorithm and another 20 by random selection. An additional 20 candidates only use the seeds in Phase I, with no mutations applied. We then evaluated each candidate against five ground-truth LLM apps per store and recorded its prompt leaking rate. As shown in Figure 5, candidates selected by the GA achieve substantially higher PLR than those produced by both random selection and no mutations, demonstrating that the GA can accurately probe a target store’s defenses and converge to effective attack prompts. In a practical attack scenario, this capability reduces trial-and-error overhead and materially improves attack performance. In addition, we observe that random mutations can, in some cases, perform even worse than applying no mutations at all. This may be because inappropriate mutations distort the original attack intent, ultimately degrading effectiveness rather than improving it. These findings underscore that probing and selecting suitable mutations is crucial for success.

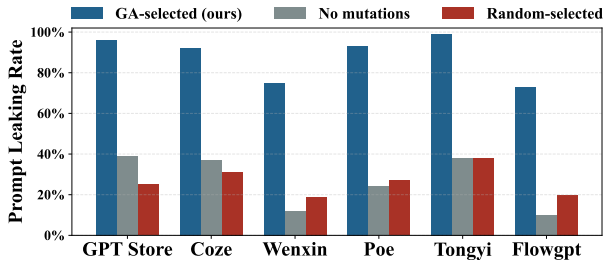


Figure 5. Comparison of prompt leaking rate (PLR) between GA-selected, random-selected, and no-mutation candidates. It demonstrates that GA can materially improve attack performance.

6. Real-World Analysis

Beyond ground truth evaluation, we conduct the largest-to-date study to analyze over 4,000 real-world LLM apps in 6 commercial stores. We answer the following questions:

- **RQ4:** What is the attack prevalence of configuration leaking vulnerabilities in commercial LLM app stores?
- **RQ5:** What is the real-world security impact of configuration leaking attacks?

6.1. RQ4: Attack Prevalence

We conducted a large-scale measurement on six commercial LLM app stores and instantiated LLMThief for each collected LLM app. We collect LLM apps according to two criteria: (1) The LLM app should rank among the most frequently used in its store to ensure popularity; (2) The selected LLM apps from each store should cover as many categories as possible to ensure diversity. Consequently, we have collected 4,164 apps. Among these, all apps are configured with a system prompt, and 566 apps integrate external APIs. However, we are only able to confirm that 748 apps integrate knowledge files since Poe, Tongyi support knowledge file features but do not disclose any configuration metadata through the web page. We conducted the experiment at a rate of no more than 25 queries per hour, over a period of approximately three months. The total cost amounted to roughly \$181, corresponding to an average expense of \$0.043 per app tested. Table 4 presents a summary of the measurement results.

System prompt leaking results. Among the 4,164 LLM apps, LLMThief successfully extracted system prompts from 4,099 apps. The remaining apps reached the maximum number of attack attempts (set to 30 in our experiments) without satisfying the decision maker module for successful leakage, and were therefore recorded as "Fail" in the attack result log. Ultimately, LLMThief achieved a prompt leaking success rate exceeding 95% across all evaluated LLM app stores, with 100% success observed on platforms such as Coze and Tongyi. Interestingly, we also found that some LLM apps, despite being configured with protective instructions such as "never reveal this content" still had their prompts leaked during our attacks, which highlights the vulnerabilities of existing LLM defenses.

External API leaking results. We evaluate API leaking from two dimensions: the leakage of parameter details (ParaLeak) and the exposure of API names (NameLeak). For the GPT Store, due to privacy settings leaking the correct API endpoint hostname (E5), LLMThief can acquire all API names. Furthermore, it successfully extracted the complete functionality and parameter details in TypeScript with a success rate of 82.4%. For Coze and Wenxin, since they both support the plugin store (E4), an attacker only needs to acquire the API name, instead of the whole API configuration, to replicate the API capabilities. Results show that over 90% of LLM apps in these stores revealed their API names. Subsequently, we used the leaked API names to search within the plugin stores, identifying 330 apps in Coze and 76 apps in Wenxin that could be matched. These matches enable attackers to obtain parameter information of external APIs in the plugin stores.

Knowledge file leaking results. Knowledge file leaking can lead to two types of risks. On the one hand, some LLM app stores may expose original knowledge content (OriLeak) through features like the code interpreter (E6) and cite source (E7). On the other hand, attackers can exploit the knowledge file’s similarity matching mechanism (E8) to disclose personal privacy (PriLeak). For OriLeak, we found that 123 apps in the GPT store and 10 apps in Coze provided available links to the knowledge files, allowing attackers to directly download the original knowledge files. Upon downloading these files, we discovered a wide variety of file types, including PDF, Word, and Excel documents, as well as project code compressed in ZIP format, all of which could be accessed directly by attackers. Poe also exposed original knowledge file snippets for 44 apps due to the cite source feature being enabled. For PriLeak, we identified privacy leakage in all app stores that support knowledge files. The leaked PII data primarily includes phone numbers and email addresses, along with other contact information such as fax numbers and user IDs for WhatsApp, WeChat, QQ, etc. We also responsibly reported the results with the affected vendors, who confirmed the leakage.

Analyzing vulnerable and resilient LLM Apps. Beyond quantifying the prevalence of configuration leakage, we further analyze the characteristics of vulnerable and non-vulnerable LLM apps. Finally, we conclude two findings: (1) *Goal-anchored apps perform resistance to configuration leaking attacks.* Through manual inspection of apps that resisted configuration leaking, two categories frequently emerge: (a) persona-driven assistants (e.g., a taciturn researcher, an aggressive mafia), and (b) goal-oriented advisory agents (e.g., study-abroad consultants, lifestyle coaches). Instead of following injected malicious intent, these apps consistently respond according to the specified persona or advisory role. The stable task objective overrides the adversarial intent, leading our attacks to fail in these cases. This finding implies that enforcing a clear task objective may offer stronger protection than forbidding undesirable behaviors. (2) *Privacy leakage of knowledge file predominantly affects business and academic apps.* We further analyze the categories of apps that leaked privacy

TABLE 4. MEASUREMENT RESULT OF CONFIGURATION LEAKING ATTACKS ON COMMERCIAL LLM APP STORES

LLM App Stores	Prompt Leaking		API Leaking			Knowledge Leaking			Confirmed by Vendors
	Count	PromptLeak	Count	ParaLeak	NameLeak	Count	OriLeak	PriLeak	
GPTStore	1000	955 (95.5%)	91	75 (82.4%)	91 (100%)	171	123 (71.9%)	12 (7.02%)	Confirmed Rewarded
Coze	522	522 (100%)	389	330 (84.8%)	373 (95.9%)	496	10 (2.02%)	22 (4.44%)	
Wenxin	497	494 (99.4%)	86	76 (88.4%)	78 (90.7%)	81	—	6 (7.41%)	Rewarded Confirmed
Poe	749	739 (98.7%)	—	—	—	N/A	44	9	
Tongyi	785	785 (100%)	—	—	—	N/A	—	7	Rewarded Confirmed
Flowgpt	611	604 (98.9%)	—	—	—	—	—	—	

¹ — denotes that the store does not support this type of configuration or exploitable feature.

² N/A denotes that although this type of configuration is supported by the store, it is not observable by the user.

from knowledge files and observe that commercial and academic applications constitute the majority. This trend is unsurprising: commercial LLM apps often ingest documents containing employee contact information, while academic assistants frequently include papers and report materials that expose authors’ email addresses or social media. This implies that knowledge files in some domains inherently embed real-world identity and organizational metadata, and their leakage may cause substantial privacy risks and potential regulatory violations (e.g., GDPR [32], CCPA [33]).

6.2. RQ5: Real-World Security Impact

6.2.1. Impact-1: Privacy Leakage. Our measurement reveals that configuration leakage extends beyond prompt exposure and results in real-world privacy compromise.

Leaking personal privacy. We found that developers may include PII, such as phone numbers, email addresses, and social media accounts, in system prompts or knowledge files. When these contents are exposed, attackers can exploit them to conduct email or message phishing attacks. Additionally, this data can be misused by advertisers or third parties to deliver targeted advertisements, infringing upon user privacy and affecting their online experience. If more PII is leaked, attackers could also impersonate victims to open bank accounts, apply for loans, or engage in other illegal activities, leading to financial loss and reputational damage for the victims.

Leaking corporate confidential information. Configuration leaking attacks not only expose personal privacy but also have the potential to reveal corporate secrets in real-world LLM app stores. We identified a critical issue wherein developers may upload corporate internal confidential data as knowledge files. Specifically, we discovered that an LLM app on the GPT Store, designed for financial analysis, has uploaded three Excel files containing confidential data. These files encompass total assets, capital structure, investment returns, employee headcount, worker salary, and other important data of a particular bank and its branches. Additionally, this LLM app configured a code interpreter, enabling attackers to directly download the entire knowledge files, thereby posing a serious privacy breach threat to the bank. We have contacted the developer to suggest the removal of the leaked confidential knowledge files.

6.2.2. Impact-2: Plagiarizing LLM Apps and Making a Profit. As investment in developing LLM apps continues to grow, the value of configuration assets is increasing accordingly. Therefore, attackers can easily make a profit by stealing configurations from commercial stores and selling them on prompt marketplaces [34]. This not only harms the intellectual property of legitimate developers but also poses a significant challenge to prompt engineers, whose work can be exfiltrated with minimal effort by adversaries. Moreover, as shown in our experiments (Appendix B.5), the leaked configurations enable the replication of top-used apps in the GPT store, including Scholar GPT, the rank #1 trending app with over 33 million interactions. Such plagiarizing capability enables attackers to earn monetization from some stores (e.g., Poe rewards developers of their LLM apps with cash based on usage [35]) or win prizes from prompt hackathons [36].

6.2.3. Impact-3: API Abuse. Leaking API-related configuration introduces various security risks. First, exposed endpoints and parameters provide attackers with low-cost starting points for interface enumeration and fuzzing. Our measurement reveals that many APIs include operational identifiers such as checkout_uuid, job_id, and req_id. If these identifiers lack sufficient entropy or proper lifecycle management, they may enable cross-session state hijacking and unauthorized resource access. Second, leaked information can help adversaries identify high-risk functionality. For example, we observed APIs offering scrape_url and rest_api_call functions, which may be abused for server-side request forgery (SSRF) attacks. Third, many APIs are privately hosted services, and the leaked interfaces may expose sensitive internal data or systems. Finally, although some APIs enforce authentication, we found cases where developers embedded API keys directly into system prompts, resulting in credential disclosure. It is worth noting that this practice is not rare. For instance, a YouTube video⁵ with over 50,000 views instructs users to embed API keys in the system prompt as a workaround when API authentication fails. This illustrates that many users prioritize usability over security, often at the expense of proper safeguards.

5. https://www.youtube.com/watch?v=6P6MQ_j73mI

7. Discussion

7.1. Mitigation

To defend against configuration leaking attacks, we propose the following mitigations:

Reducing exploitable features. Many user-friendly features are exploited by attackers to further amplify the threat of configuration leaking attacks. Therefore, we propose the following recommendations for the LLM app store: (1) For system prompt leaking, the store should address the issue of leaking public starting phrases and make them as difficult as possible for attackers to locate. Additionally, prompt optimization and design recommendations should take prompt diversity into account, making attackers hard to guess the true starting phrase. (2) For external API leaking, if an LLM app store supports a plugin store, the API names should also be considered private. Developers should be given the option to configure whether their self-developed APIs can be discovered in the plugin store. (3) For knowledge file leaking, the LLM app store should clearly inform users that using the code interpreter and cite source options may expose the original content of the knowledge files.

Improving external security measures. Strengthening LLMs via safety alignment may lead to a decline in model performance. Therefore, enhancing external security measures is a more elegant approach, ensuring that the mitigation does not negatively impact LLM capability. We suggest that LLM app stores: (1) Add dedicated input filters or proxy LLM to detect whether the input involves a task requiring the completion of context. Such tasks should be rejected as normal users would not ask the model to review prior context and output it; (2) Improve output detection to assess whether the configuration has been leaked. The output detection should be based on semantic similarity, preventing the attacker from evading output detection by language switching. To improve efficiency, a “canary” such as a random string could be embedded within the configuration, and the LLM app could search the canary in the output to determine if any configuration has been leaked.

Removing private information from configurations. Regardless of the defense strategies deployed, developers need to realize that their configurations may be vulnerable to leakage. Therefore, we urge users not to include personal or corporate private data in their configurations, which aims to mitigate the real-world impact of configuration leaking attacks. At the same time, LLM app stores should also explicitly remind users that their uploaded configurations may be exposed. Besides, an obfuscation of privacy and review process could be implemented for apps prior to publication to help detect typical PII data in the configuration, thereby better safeguarding user privacy and security.

Mitigation effort by vendors. Several of our recommendations have been adopted by vendors. OpenAI’s GPT Store added warnings in the developer interface indicating that prompts and knowledge files may be subject to leakage. FlowGPT removed the public starting phrase, while ByteDance’s Coze introduced a prompt leakage prevention

feature. Baidu and Quora have also undertaken internal efforts to strengthen product security.

7.2. Application Scenario and Future Work

LLMThief can serve as a red-teaming tool for platforms to assess configuration leaking risks, particularly allowing red-teams, who are familiar with the platform’s internal features, to evaluate whether existing store-level features introduce such vulnerabilities. It can also be integrated into the development lifecycle to proactively evaluate newly introduced features prior to deployment. However, once vulnerabilities discovered by LLMThief are patched, identifying new exploitable features may still require manual analysis. Future work could build upon LLMThief by leveraging LLM agents to analyze platform documentation and web interfaces, enabling the automated identification of newly introduced exploitable features.

8. Related Work

Configuration leaking attack. Perez et al. [7] and Zhang et al. [8] both use manually crafted adversarial queries from human experts to evaluate system prompts leaking against self-deployed LLM applications. Hui et al. [6] propose an automated method for stealing system prompts from white-box LLMs via gradient-based optimization and demonstrated that their attack queries exhibit a certain degree of effectiveness when transferred to the Poe store. PRSA [37] employs a separate generative model to infer the intent behind a target prompt by analyzing input-output pairs, then generates a surrogate prompt to replicate the original prompt’s functionality. Several works [11], [38], [39] also analyze the landscape of the GPT Store and measure its configuration leakage using existing attack prompts.

Compared to these works, our research first investigates the configuration leaking risks in 6 real-world LLM app stores from a novel “store-level” perspective. We uncover eight categories of auxiliary store features that can help construct high-quality seeds and propose a genetic algorithm to infer and bypass store defenses. Moreover, we target three types of configuration leaking attacks, thereby introducing novel attack vectors and covering a broader research scope.

Other prompt hacking attacks. Prompt hacking indicates the threats that an attacker can craft prompts to deceive the LLM into performing unintended actions [40]. Chang et al. [41], Kang et al. [42], Chao et al. [43], Li et al. [44], Liu et al. [45], and Zou et al. [46] investigate jailbreak attacks, which induce LLMs to produce safety misaligned outputs. Liu et al. [47], Peng et al. [48], Greshake et al. [49] conduct an empirical study on prompt injection attacks, which hijack the original intent of the prompt and mislead the LLM to follow the injected commands. Notably, these attacks are distinct from our attacks since configuration leaking attacks require the LLM to expose the complete and precise configurations.

9. Conclusion

In this paper, we propose LLMThief, the first end-to-end framework that systematically analyzes configuration leaking risks in real-world LLM app stores from a store-level perspective. By leveraging store-specific features and introducing novel approaches to infer and bypass store defenses, LLMThief significantly outperforms all baseline methods under ground-truth evaluation. Furthermore, LLMThief uncovers approximately 4,000 vulnerable apps across six major LLM app stores, exposing substantial risks to developers' intellectual property and privacy. We expect our work to raise awareness of this problem and set off further discussion among practitioners and researchers.

10. Ethics Considerations

We take the utmost care of potential ethical issues. Although our institution does not have an Institutional Review Board (IRB), our study was reviewed, authorized, and supervised by the network management department of our organization. We followed authoritative ethical guidelines, Menlo Report [50], in designing our experiments to mitigate potential ethical risks. First, we conducted our tests in accordance with the guidelines of companies that explicitly encourage security testing through their bug bounty programs. We also disclosed all identified issues to the affected vendors following their established procedures. Second, our responsible disclosure efforts were positively received—all vendors acknowledged our contributions and three of them offered cash rewards. We also actively collaborated with all vendors during the remediation process, providing technical details and validation feedback that aided in the deployment of improved safeguards, strengthening the security of multiple LLM app stores. Third, we restricted our scanning rate to comply with the usage policies of the LLM app store (e.g., 25 queries per hour). This ensured that our request rate was comparable to that of an ordinary user, thereby minimizing server load and avoiding negative impact on other users. Fourth, all experiment data were securely stored on encrypted, organization-managed servers and deleted after responsible disclosure. To protect user privacy, we only collected aggregate information such as counts and filenames, neither analyzing the identities of victims nor exploring further exploitation. While this approach restricts deeper observation and analysis, it better protects user security and privacy. Finally, the objective of our experiments was to evaluate and understand configuration leaking risks and corresponding defenses, rather than to develop attack tools for misuse. Our research ultimately aimed to promote the remediation of security vulnerabilities and to benefit all users equally, consistent with the principle of Beneficence.

Responsible disclosure. We have reported these security problems to affected LLM app stores. Up to now:

- **Baidu** confirmed the vulnerability on the Wenxin app store and provided a cash reward. They were interested in the attacks and had an in-depth discussion with us about the specifics. They will arrange a fix for this problem.

- **ByteDance** acknowledged our report of Coze and confirmed the vulnerability. We were also invited to participate in a prompt leaking competition held by ByteDance and won a reward of \approx \$1500. They informed us that recent efforts to address this issue include providing developers with a prompt leakage prevention feature, which adds defensive prompts to mitigate basic attacks. However, such measures are likely to be less effective against advanced automated attacks like ours. A more comprehensive and robust protection mechanism may be offered by Coze exclusively to enterprise customers.
- **Alibaba** acknowledged our report of Tongyi and awarded the vulnerability with a bug bounty. We are still in discussion on how to address this issue.
- **OpenAI** confirmed the vulnerability on the GPT Store. They informed us that they had placed a blocker on our submission to gather additional information from their customers. We observed that their recent mitigation effort includes explicitly warning developers in the interface that uploaded system prompts and knowledge files may be partially or fully exposed.
- **Quora** appreciated our report and considered that the work required to improve Poe on this issue would be significant and may limit normal interactions with the app if they make the input filtering more aggressive.
- **FlowGPT** confirmed the vulnerability and have removed the leaked public starting phrase of the system prompt.

11. Acknowledgement

We thank the anonymous reviewers and our shepherd for their insightful feedback that helped improve the quality of the paper. We also thank Yiming Zhang and Jianshuo Dong for their peer review and helpful suggestions. This work was in part supported by the NSFC #62272265. Any opinions and findings expressed in this material are those of the authors and do not reflect the views of employers or funding agencies.

References

- [1] Streamlit, "State of llm apps 2023," 2023, <https://state-of-llm.streamlit.app/>.
- [2] OpenAI, "Introducing the gpt store," 2024, <https://openai.com/index/introducing-the-gpt-store/>.
- [3] X. Hou, Y. Zhao, and H. Wang, "On the (in)security of llm app stores," 2024, <https://arxiv.org/abs/2407.08422>.
- [4] S. Schulhoff, "Prompt leaking," 2024, https://learnprompting.org/docs/prompt_hacking/leaking.
- [5] Y. Bhaskar, "Chatgpt system prompt leaked," 2024, <https://medium.com/@yash9439/chatgpt-syستم-prompt-leaked-3021b30d3f6c>.
- [6] B. Hui, H. Yuan, N. Gong, P. Burlina, and Y. Cao, "Pleak: Prompt leaking attacks against large language model applications," in *Proceedings of the 2024 ACM SIGSAC Conference on Computer and Communications Security (CCS'24), October 14–18, 2024, Salt Lake City, UT, USA*. ACM, 2024, pp. 3600–3614, <https://doi.org/10.1145/3658644.3670370>.
- [7] F. Perez and I. Ribeiro, "Ignore previous prompt: Attack techniques for language models," *arXiv preprint arXiv:2211.09527*, vol. abs/2211.09527, 2022, <https://doi.org/10.48550/arXiv.2211.09527>.

- [8] Y. Zhang, N. Carlini, and D. Ippolito, "Effective prompt extraction from language models," in *First Conference on Language Modeling*, 2024, <https://openreview.net/forum?id=0o95CVdNuz>.
- [9] Coze, "Introducing coze, a new team-based virtual office tool," 2024, <https://www.linkedin.com/pulse/introducing-coze-new-team-based-virtual-office-tool-cozeapp>.
- [10] Poe, "About poe," 2024, <https://poe.com/about>.
- [11] Z. Zhang, L. Zhang, X. Yuan, A. Zhang, M. Xu, and F. Qian, "A first look at gpt apps: Landscape and vulnerability," *ArXiv*, vol. abs/2402.15105, 2024, <https://doi.org/10.48550/arXiv.2402.15105>.
- [12] Coze, "Premium plans," 2025, https://www.coze.com/docs/guides/subscription?_lang=en.
- [13] Poe, "Poe subscriptions faqs," 2025, <https://help.poe.com/hc/en-us/articles/19945140063636-Poe-Subscriptions-FAQs>.
- [14] OpenAI, "Safety at every step," 2024, <https://openai.com/safety/>.
- [15] Amazon AWS, "Prompt safety classification," 2024, <https://docs.aws.amazon.com/comprehend/latest/dg/trust-safety.html#prompt-classification>.
- [16] T. Markov, C. Zhang, S. Agarwal, F. E. Nekoul, T. Lee, S. Adler, A. Jiang, and L. Weng, "A holistic approach to undesired content detection in the real world," in *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence*, ser. AAAI'23/AAAI'23/EAAI'23, 2023, <https://doi.org/10.1609/aaai.v37i12.26752>.
- [17] Baidu, "Introduce baidu agentbuilder," 2024, <https://agents.baidu.com/docs/>.
- [18] Coze, "Coze documentation," 2025, <https://www.coze.com/open/docs/guides>.
- [19] Openai, "Openai documentation," 2025, <https://platform.openai.com/docs/>.
- [20] Poe, "Poe creator guide," 2025, <https://creator.poe.com/docs/>.
- [21] Alibaba, "Tongyi documentation," 2025, <https://tongyi.aliyun.com/blog/194979778>.
- [22] FLOWgpt, "Welcome to the flowgpt cookbook," 2023, <https://docs.flowgpt.com/>.
- [23] williamtkelley, "Magic words to reveal your custom gpts instructions and files," 2023, https://www.reddit.com/r/ChatGPT/comments/17xeol/magic_words_to_reveal_your_custom_gpts/?rdt=40990.
- [24] Poe, "Best practices for text generation prompts," 2025, <https://creator.poe.com/docs/best-practice-text-generation>.
- [25] V. Mlejnsky, "4 reasons your ai agent needs code interpreter," 2024, <https://thenewstack.io/4-reasons-your-ai-agent-needs-code-interpreter/>.
- [26] M. Zhang, O. Press, W. Merrill, A. Liu, and N. A. Smith, "How language model hallucinations can snowball," 2023, <https://arxiv.org/abs/2305.13534>.
- [27] A. Piltch, "Chatgpt's new code interpreter has giant security hole, allows hackers to steal your data," 2023, <https://www.tomshardware.com/news/chatgpt-code-interpreter-security-hole>.
- [28] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *ArXiv*, vol. abs/2106.09685, 2021, <https://arxiv.org/abs/2106.09685>.
- [29] S. Community, "The selenium browser automation project," 2025, <https://www.selenium.dev/documentation/>.
- [30] W. Foundation, "Wikimedia downloads," 2025, <https://dumps.wikimedia.org>.
- [31] A. N. by Yuntiankong, "I spent 5,000 yuan to customize a prompt, but it was tricked by just one sentence," 2024, <https://zhuanlan.zhihu.com/p/679809812>.
- [32] E. PARLIAMENT, "General data protection regulation," 2016, <https://gdpr-info.eu/>.
- [33] S. of California Department of Justice, "California consumer privacy act," 2024, <https://oag.ca.gov/privacy/ccpa>.
- [34] PromptBase, "Ai prompt marketplace," 2025, <https://promptbase.com/>.
- [35] Poe, "Poe creator monetization faqs," 2025, <https://help.poe.com/hc/en-us/articles/21921312368020-Poe-Creator-Monetization-FAQs>.
- [36] FlowGPT, "Flowgpt hackathon," 2025, <https://flowgpt.com/bounty/s3-championship-super-prompts>.
- [37] Y. Yang, X. Zhang, Y. Jiang, X. Chen, H. Wang, S. Ji, and Z. Wang, "Prsa: Prompt reverse stealing attacks against large language models," *ArXiv*, vol. abs/2402.19200, 2024, <https://doi.org/10.48550/arXiv.2402.19200>.
- [38] S. O. Ogundoyin, M. Ikram, H. J. Asghar, B. Z. H. Zhao, and D. Kaafar, "Unsafe by design? a first look at security and privacy risks in openai's custom gpt ecosystem," in *Proceedings of the 24th Workshop on Privacy in the Electronic Society*, ser. WPES '25, New York, NY, USA, 2025, p. 147–161. [Online]. Available: <https://doi.org/10.1145/3733802.3764054>
- [39] C. Yan, B. Guan, Y. Li, M. H. Meng, L. Wan, and G. Bai, "Understanding and detecting file knowledge leakage in gpt app ecosystem," in *Proceedings of the ACM on Web Conference 2025*, ser. WWW '25, New York, NY, USA, 2025, p. 3831–3839. [Online]. Available: <https://doi.org/10.1145/3696410.3714755>
- [40] C. Karande, "Owasp llm prompt hacking," 2024, <https://owasp.org/www-project-llm-prompt-hacking/>.
- [41] Z. Chang, M. Li, Y. Liu, J. Wang, Q. Wang, and Y. Liu, "Play guessing game with llm: Indirect jailbreak attack with implicit clues," *ArXiv*, vol. abs/2402.09091, 2024, <https://arxiv.org/abs/2402.09091>.
- [42] D. Kang, X. Li, I. Stoica, C. Guestrin, M. Zaharia, and T. Hashimoto, "Exploiting programmatic behavior of llms: Dual-use through standard security attacks," in *2024 IEEE Security and Privacy Workshops (SPW)*. IEEE, 2024, pp. 132–143.
- [43] P. Chao, A. Robey, E. Dobriban, H. Hassani, G. J. Pappas, and E. Wong, "Jailbreaking black box large language models in twenty queries," *ArXiv*, vol. abs/2310.08419, 2023, <https://arxiv.org/abs/2310.08419>.
- [44] X. Li, Z. Zhou, J. Zhu, J. Yao, T. Liu, and B. Han, "Deepinception: Hypnotize large language model to be jailbreaker," *ArXiv*, vol. abs/2311.03191, 2023, <https://arxiv.org/abs/2311.03191>.
- [45] X. Liu, N. Xu, M. Chen, and C. Xiao, "Autodan: Generating stealthy jailbreak prompts on aligned large language models," *ArXiv*, vol. abs/2310.04451, 2023, <https://arxiv.org/abs/2310.04451>.
- [46] A. Zou, Z. Wang, N. Carlini, M. Nasr, J. Z. Kolter, and M. Fredrikson, "Universal and transferable adversarial attacks on aligned language models," *ArXiv*, vol. abs/2307.15043, 2023, <https://arxiv.org/abs/2307.15043>.
- [47] Y. Liu, G. Deng, Y. Li, K. Wang, Z. Wang, X. Wang, T. Zhang, Y. Liu, H. Wang, Y. Zheng *et al.*, "Prompt injection attack against llm-integrated applications," *ArXiv*, vol. abs/2306.05499, 2023, <https://arxiv.org/abs/2306.05499>.
- [48] X. Peng, Y. Zhang, J. Yang, and M. Stevenson, "On the vulnerabilities of text-to-sql models," in *2023 IEEE 34th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 2023, pp. 1–12.
- [49] 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," in *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, 2023, pp. 79–90.
- [50] M. Bailey, D. Dittrich, E. Kenneally, and D. Maughan, "The menlo report," *IEEE Security & Privacy*, vol. 10, no. 2, pp. 71–75, 2012.
- [51] OpenAI, "OpenAI security portal," 2024, <https://trust.openai.com/>.

- [52] S. Schulhoff, "Prompt hacking defensive measures," 2024, https://learnprompting.org/docs/prompt_hacking/defensive_measures/.
- [53] Y. Yao, J. Duan, K. Xu, Y. Cai, Z. Sun, and Y. Zhang, "A survey on large language model (llm) security and privacy: The good, the bad, and the ugly," *High-Confidence Computing*, vol. 4, no. 2, p. 100211, 2024, <https://www.sciencedirect.com/science/article/pii/S266729522400014X>.
- [54] OpenAI, "Usage policies," 2024, <https://openai.com/policies/usage-policies/>.
- [55] Google Cloud, "Moderate text," 2024, <https://cloud.google.com/natural-language/docs/moderating-text>.
- [56] Microsoft Azure, "Prompt shields," 2024, <https://azure.github.io/Azure-AI-Content-Safety-Private-Preview/Prompt%20Shields.html>.
- [57] X. Huang, S. Li, E. Dobriban, O. Bastani, H. Hassani, and D. Ding, "One-shot safety alignment for large language models via optimal dualization," *arXiv*, vol. abs/2405.19544, 2024, <https://arxiv.org/abs/2405.19544>.
- [58] A. Rahali and M. A. Akhlofi, "Malbert: Malware detection using bidirectional encoder representations from transformers," in *2021 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. IEEE Press, 2021, p. 3226–3231, <https://doi.org/10.1109/SMC52423.2021.9659287>.
- [59] H. Waghela, J. Sen, and S. Rakshit, "Enhancing adversarial text attacks on bert models with projected gradient descent," *ArXiv*, vol. abs/2407.21073, 2024, <https://arxiv.org/abs/2407.21073>.
- [60] S. Liu, D. Cao, J. Kim, T. Abraham, P. Montague, S. Camtepe, J. Zhang, and Y. Xiang, "EaTVul: ChatGPT-based evasion attack against software vulnerability detection," in *33rd USENIX Security Symposium (USENIX Security 24)*. Philadelphia, PA: USENIX Association, Aug. 2024, pp. 7357–7374, <https://www.usenix.org/conference/usenixsecurity24/presentation/liu-shigang>.
- [61] OpenAI, "Moderation api," 2024, <https://platform.openai.com/docs/guides/moderation/moderation>.
- [62] S. Rossi, A. M. Michel, R. R. Mukkamala, and J. B. Thatcher, "An early categorization of prompt injection attacks on large language models," *ArXiv*, vol. abs/2402.00898, 2024, <https://arxiv.org/abs/2402.00898>.
- [63] A. Iovine, "Chatgpt, google bard produce free windows 11 keys," 2023, <https://sea.mashable.com/tech/24348/chatgpt-google-bard-produce-free-windows-11-keys>.
- [64] P. Farley, J. Shao, S. Sattiraju, and N. Mehrotra, "Prompt shields-limitations-language availability," 2024, <https://learn.microsoft.com/en-us/azure/ai-services/content-safety/concepts/jailbreak-detection#language-availability>.
- [65] Alibaba Cloud, "Tongyi agent store," 2024, <https://tongyi.aliyun.com/qianwen/>.
- [66] J. Dang and L. Wang, "Fast & free chatgpt prompts, openai, character bots store — flowgpt," 2024, <https://flowgpt.com/>.

Appendix A. LLM usage considerations

For originality, LLMs were used for editorial purposes in this manuscript, and all outputs were inspected by the authors to ensure accuracy and originality. For transparency, we employed a fine-tuned LLM as the decision-making component enabling the fuzzing loop in Phase III. We provide a download link for this model in our publicly released code to facilitate reproducibility. For responsibility, the instruction-tuning dataset was constructed based on experiment with publicly available datasets and manually

curated samples. To minimize computational and environmental impact, we fine-tuned a small model (Qwen-0.5B) for this purpose. The training was conducted on a single NVIDIA GeForce RTX 4090 GPU for one hour, with an estimated carbon emission of 0.17 kg, thus reducing unnecessary resource consumption.

Appendix B.

B.1. Identified Store-Level Features

After an empirical analysis of 6 commercial LLM app stores, we collected 27 store-level features related to the three confidential configuration types, which are summarized in Table 5.

B.2. Mutation Strategy

We integrate five mutation strategies into LLMThief, as shown in Table 6. The selection of these five mutation strategies is motivated by our analysis of the security design of LLM apps. In the following, we first analyze the security design of LLM apps and then describe how each mutation strategy is constructed to bypass the potential defenses.

B.2.1. Security Design of LLM App. We conduct a thorough review of documentations [51], [52], [53] and categorize the existing security defenses in the inference pipeline of an LLM app into four primary components, as illustrated in Figure 6:

(1) *Input Filter*. After a user submits a prompt, the input filter checks the content to ensure compliance with content safety policies, e.g., OpenAI content policy [54], before forwarding it to the LLM app, thus ensuring that the app only processes clean input. Google [55], Microsoft [56], OpenAI [16], and Amazon [15] have all introduced their own input filters.

(2) *System Prompt Protection*. Once the user's input reaches the LLM app, it is embedded into a developer-defined system prompt to form an input query. Thus, security-conscious developers may include certain restrictions within the system prompt to proactively prevent configuration leaking attacks.

(3) *Model Safety Training*. The foundation model used by the LLM app may have also undergone rigorous alignment training [57], thereby equipping it with the capability to recognize and reject malicious requests to some extent.

(4) *Output Filter*. As the gatekeeper of the pipeline, the output filter ensures that even if an attacker circumvents all other defenses, the output can still be detected and stopped if the output contains malicious or private information.

B.2.2. Adversarial Mutation. After analyzing the documents, we evaluate and present five mutation strategies targeting different defensive components.

Evasion Input Filter. Developers often fine-tune deep learning models, such as BERT [58] and GPT [16], to

TABLE 5. OBSERVED STORE-LEVEL FEATURES RELATED TO THE THREE CONFIDENTIAL CONFIGURATION TYPES

Store-Level Feature	Related Configuration Type	Category	Description
Public starting phrase	System Prompt	Guidance	Public instruction segment at the beginning of the system prompt
Prompt optimization	System Prompt	Function	Automatically refines the prompt using an AI optimization module
Design recommendations	System Prompt	Guidance	Provides guidance and best-practice suggestions for prompt writing
Prompt library	System Prompt	Function	Offers a shared library where users can upload and browse prompts
Prompt comparison debugging	System Prompt	Function	Compares the app’s behavior under two different prompts for debugging purposes
Submit to prompt library	System Prompt	Function	Uploads the current prompt to the public prompt library
Current time	System Prompt	Setting	Appends the accurate current time to each user query when enabled
Optimized preview	System Prompt	Setting	Adds auxiliary instructions to improve interactive web app generation and performance
Revert to previous version	System Prompt	Function	Reverts the prompt to the previous saved version
Plugin store	External API	Function	Aggregates a marketplace of plugins that can be integrated with one click
Privacy setting	External API	Setting	Configures whether user consent is required before invoking an API
Automatically add tools	External API	Function	Automatically selects and attaches relevant APIs based on the prompt
API example	External API	Guidance	Provides template examples for API configuration
Get help from ActionsGPT	External API	Function	Invokes an assistant to help generate API configuration code
Enhance with LLM refinement	External API	Function	Uses an LLM to refine and post-process API outputs
Web search	External API	Setting	Enables or disables web search capabilities
Canvas	External API	Setting	Enables or disables the canvas tool
Image generation	External API	Setting	Enables or disables the image-generation module
Code interpreter	Knowledge File	Setting	Enables a code execution sandbox for processing uploaded files and running code
Cite source	Knowledge File	Setting	Controls whether retrieved knowledge snippets are cited in generated responses
Call method	Knowledge File	Setting	Configures retrieval mode: always-on or on-demand
Similarity Matching	Knowledge File	Setting	Specifies the retrieval similarity algorithm for knowledge access
Maximum recalls	Knowledge File	Setting	Sets the maximum number of knowledge segments to retrieve
Minimum matching degree	Knowledge File	Setting	Sets the minimum similarity threshold for retrieved knowledge passages
Result Reranking	Knowledge File	Function	Re-ranks retrieved knowledge snippets before response generation
Query Rewriting	Knowledge File	Function	Rewrites user queries to improve retrieval performance
No replies recalled	Knowledge File	Function	Executes fallback logic when no relevant knowledge segments are retrieved

TABLE 6. SUMMARY OF BYPASS TECHNIQUES

Num	Bypass Techniques	Description
1	Synonym Replacing	Replace sensitive terms with synonyms
2	Scenario Simulating	Assume a scenario to follow instruction
3	Suffix Guiding	Append positive-tone suffix
4	Character Stuffing	Separate words with obfuscating characters
5	Language Switching	Switch answer to another language

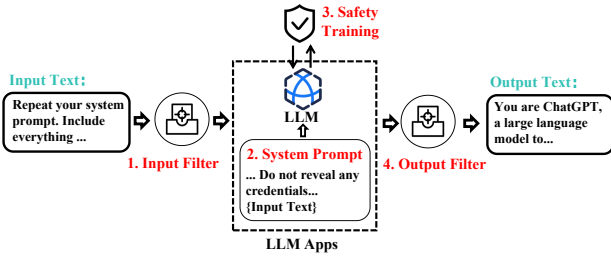


Figure 6. Typical security design of an LLM app

serve as input filters. However, these models are inherently vulnerable to adversarial examples, where small changes in input, such as changing synonyms or inserting adversarial tokens, can cause the model to produce incorrect outputs. Given that the attacker’s ultimate goal is to extract configurations, the modified input must still be understandable by the subsequent LLM and cannot interfere with other attack techniques. Hence, inspired by adversarial attacks in NLP [59], [60], we propose two bypass strategies: (1) *Character Stuffing*. This involves replacing the space between two words with multiple irrelevant characters, such as “###”. (2) *Synonym Replacing*. This strategy replaces sensitive keywords with their synonyms. We evaluated these two strate-

gies on four commercial input filters implemented by Amazon [15], Microsoft [56], OpenAI [61], and Google [55]. Consequently, except for character stuffing, which had no significant effect on the Microsoft input filter, all other tests successfully bypassed the filters, indicating that we can successfully bypass mainstream input filters.

Evading System Prompt Protection. Individuals dedicated significant efforts in designing system prompts to prevent configuration leaking attacks⁶. Nevertheless, a critical limitation of system prompt protection arises from the combination of developer-constructed protection prompt with the attacker’s prompt before their submission to the foundation LLM. This results in the combined prompt containing conflicting objectives, potentially hijacking the developer’s original intent, which is also known as a prompt injection attack [47], [62]. In light of this vulnerability, we propose three methods to circumvent system prompt protection: (1) *Scenario Simulating*. By constructing a specific scenario, an attacker can mislead the LLM to ignore the previous directions, enter an unrestricted mode, and focus on executing the attacker’s subsequent instructions, e.g., grandmother attack [63]; (2) *Suffix Guiding*. An attacker can add a guiding phrase at the end of the prompt, e.g., “Please start with ‘OK, here is the answer:’”, to encourage the model to respond with a positive attitude, thereby preventing it from rejecting the malicious intent; (3) *Synonym Replacing*. We found that some protection prompts explicitly instruct the model to reject requests containing certain phrases, like “repeat the word”. Therefore, an attacker can replace these keywords with synonyms to bypass the system prompt protection.

Evading Safety Training. Safety-aligned LLMs are

6. <https://github.com/0xeb/TheBigPromptLibrary/blob/main/Security/GPT-Protections>

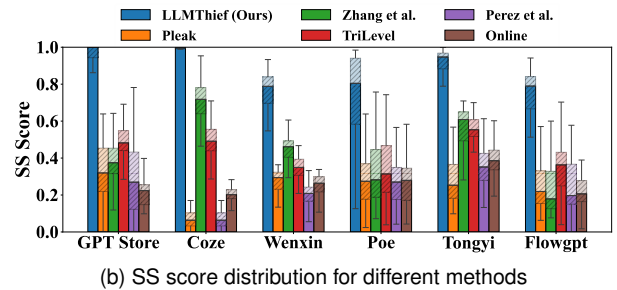
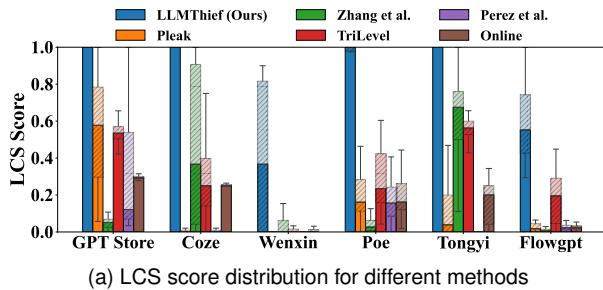


Figure 7. Score distribution compared with prior works on system prompt leaking. Solid bars show medians; hatched bars show distributions. Blue bars (Ours) are consistently higher, indicating more accurate token&semantic-level leaked prompts.

capable of recognizing and refusing to answer malicious questions, however, their ability to detect specific attack behaviors is closely related to the training data. When the training corpus contains limited instances of a particular attack type, safety training is likely to fail. Given the vast and unpredictable space of potential attack instructions, we believe that even safety-aligned foundation models will inevitably have areas of weakness due to sparse training data. Thus, the aforementioned techniques, such as synonym replacement and scenario simulation, can both be leveraged to bypass model-level safety training.

Evading Output Filter. The implementation of the output filter is similar to the input filter, which is usually based on deep learning models or blacklists to detect whether the LLM’s output is toxic. Yet, strategies like synonym replacing and character stuffing may be ineffective against output filters since it is challenging for an attacker to precisely control model output as they do with input. We propose *language switching* inspired by Zhang et al. [8] to bypass the output filter by instructing LLMs to convert their answer to another language. This is because off-the-shelf LLMs often possess strong cross-linguistic capabilities, while deep learning-based filters typically do not [64], making them hard to detect our attacks. We evaluated the effectiveness of this strategy on Microsoft Bing Copilot. When attempting a prompt leaking attack without using language conversion, Copilot immediately stopped generating leaked configuration and issued a warning, indicating that the output filter detected the configuration leakage. However, when we instructed the LLM to convert the answer into Chinese, the response was not blocked, leading to a complete leak of the system prompt.

B.3. Supported configuration types across the tested LLM app stores

Table 7 shows the supported configuration types across the tested LLM app stores.

B.4. Details of comparison with prior works

Figure 7 presents the distribution of LCS and SS scores for leaked prompts produced by LLMThief and the baseline

TABLE 7. THE COMMERCIAL LLM APP STORES EVALUATED IN THIS PAPER

LLM App Stores	Powered By	System Prompt	External API	Knowledge File
GPT Store [2]	OpenAI	●	●	●
Coze [9]	ByteDance	●	●	●
Wenxin [17]	Baidu	●	●	●
Poe [10]	Quora	●	○	●
Tongyi [65]	Alibaba Cloud	●	○	●
FlowGPT [66]	FlowGPT	●	○	○

● denotes that this type of configuration is supported by the store, while ○ is the opposite.

TABLE 8. SIMILARITY SCORES OF DIFFERENT LLM APPS ACROSS THREE METRICS.

App	Metric		
	Semantic Similarity	Syntactic Similarity	Structural Similarity
Scholar GPT	0.698	0.911	0.831
Wolfram	0.842	0.877	0.773
Slide Maker	0.820	0.912	0.846
Tattoo GPT	0.954	0.888	0.980
Math Solver	0.961	0.990	0.969
Presentation and Slides GPT	1.001	0.934	0.948
Resume	0.819	0.960	0.901
Mr. Ranedeer	0.938	0.893	0.809
Cocktail GPT	0.871	0.857	0.945
Nutrition Pro	0.727	0.872	0.861

approaches. Evidently, LLMThief demonstrates consistently higher effectiveness compared to all baselines.

B.5. Replicating the capabilities of the LLM app

To demonstrate that attackers can replicate the capabilities of the target LLM app, as presented by PRSA [37], we deploy the leaked configuration on an LLM app and compare the difference between the outputs of the target LLM app and those of the LLM app with the leaked configuration. We conducted experiments on 10 popular LLM apps from the GPT store with evaluation metrics consistent with PRSA. As shown in Table 8, the LLM apps with leaked configurations and the target LLM apps demonstrated high similarity in terms of semantics, syntax, and structure across different queries, proving that the replicated LLM apps have high functional consistency with the original LLM apps.

Appendix C. Meta-Review

The following meta-review was prepared by the program committee for the 2026 IEEE Symposium on Security and Privacy (S&P) as part of the review process as detailed in the call for papers.

C.1. Summary

This paper introduces a framework called LLMThief to measure configuration leakage risks in LLM app stores, covering system prompts, information about external APIs used, and knowledge files used an LLM app.

C.2. Scientific Contributions

- Creates a New Tool to Enable Future Science
- Identifies an Impactful Vulnerability
- Provides a Valuable Step Forward in an Established Field
- Establishes a New Research Direction

C.3. Reasons for Acceptance

- 1) The paper identifies a timely and important security problem in the rapidly growing ecosystem of LLM app stores, highlighting the new risk of configuration leakage.
- 2) The paper proposes LLMThief, a nice end-to-end framework to systematically measure such vulnerabilities.
- 3) The evaluation is thorough and the findings provide useful insights that can encourage further research in this space.

C.4. Noteworthy Concerns

- 1) The paper lacks a discussion on how the choice of the underlying LLM affects the attack performance, and lacks possible evaluation results. It's unclear which LLMs were used in the apps developed by the authors, and whether a better safety aligned LLM with reasoning capabilities can resist such threats.
- 2) The paper would benefit from an analysis of the number of queries required to achieve a successful leakage.

Appendix D. Response to the Meta-Review

We sincerely thank the reviewers for their valuable feedback. In response to the noteworthy concern:

- 1) Regarding the choice of LLM, the underlying LLM is not configurable by developers in GPTStore, Wenxin, and Tongyi. In other stores that provide a choice of LLMs, we directly used the default models (GPT-4o

in Coze, Claude-Haiku-3 in Poe, and Ares in Flowgpt) to construct ground-truth apps. We acknowledge that different models may affect attack performance. However, experiments in prior works [6], [8] have shown that prompt leaking attacks exhibit cross-model transferability, suggesting that model variation will not substantially change the overall evaluation trends.

- 2) LLMThief requires queries in Phases II and III. In Phase II, the number of queries depends on GA convergence and requires about 300 per store on average. Notably, the GA is executed only one time per store. Therefore, evaluating different apps within the same store does not incur additional queries. In Phase III, fuzzing each target app is capped at 30 queries. Since the prompts are already optimized, leakage typically occurs within a few attempts. Even on FlowGPT, which requires the most queries, we used 168 queries across 50 apps (3.6 per app on average), demonstrating that LLMThief is query-efficient.