# ObliInjection: Order-Oblivious Prompt Injection Attack to LLM Agents with Multi-source Data

Reachal Wang
Duke University
reachal.wang@duke.edu

Yuqi Jia
Duke University
yuqi.jia@duke.edu

Neil Zhenqiang Gong
Duke University
neil.gong@duke.edu

*Abstract*—Prompt injection attacks aim to contaminate the input data of an LLM to mislead it into completing an attacker-chosen task instead of the intended task. In many applications and agents, the input data originates from multiple sources, with each source contributing a segment of the overall input. In these multi-source scenarios, an attacker may control only a subset of the sources and contaminate the corresponding segments, but typically does not know the order in which the segments are arranged within the input. Existing prompt injection attacks either assume that the entire input data comes from a single source under the attacker's control or ignore the uncertainty in the ordering of segments from different sources. As a result, their success is limited in domains involving multi-source data.

In this work, we propose *ObliInjection*, the first prompt injection attack targeting LLM applications and agents with multi-source input data. ObliInjection introduces two key technical innovations: the *order-oblivious loss*, which quantifies the likelihood that the LLM will complete the attacker-chosen task regardless of how the clean and contaminated segments are ordered; and the *orderGCG algorithm*, which is tailored to minimize the order-oblivious loss and optimize the contaminated segments. Comprehensive experiments across three datasets spanning diverse application domains and twelve LLMs demonstrate that ObliInjection is highly effective, even when only one out of 6–100 segments in the input data is contaminated. Our code and data are available at: https://github.com/ReachalWang/ObliInjection.

## I. INTRODUCTION

An LLM takes a prompt as input and produces a response. The prompt typically consists of an *instruction* and a *data sample*. In many application and agent scenarios, this data sample originates from multiple sources, which we refer to as *multi-source data*. Each portion of data from a source is called a *segment*, and the data is a concatenation of segments from different sources. For example, in review summarization adopted by Amazon [1], the instruction might be: "Please summarize the following reviews:", with a product's reviews themselves forming the data. In this case, each segment corresponds to an individual review. For AI Overviews in search that summarize a news event [2], the data could consist of news articles from various outlets covering the same event, with

each segment representing one article. Similarly, in retrieval-augmented generation (RAG) systems, the data consists of passages retrieved from a knowledge database, with each retrieved passage treated as a segment. When an LLM agent selects a tool (e.g., an MCP server) from the available options, the data includes the user's task description as well as the names and descriptions of the available tools, with each tool's name–description pair forming a segment.

Due to the inseparability of instructions and data in a prompt, combined with the strong instruction-following capabilities of LLMs, these models are fundamentally vulnerable to *prompt injection attacks* [3]–[6]. Specifically, when the data originates from untrusted sources, an attacker can embed a malicious prompt into it, causing the LLM to produce an attacker-chosen response that completes an attacker-chosen task rather than the intended task. We refer to the attacker-chosen task as the *injected task* and the intended task as the *target task*. For example, in review summarization, the attacker-chosen response could be "The product is useless!" misleading the LLM to generate a summary that could damage the product's reputation. Major technology companies [7]–[10] now routinely conduct extensive vulnerability testing against prompt injection attacks before releasing or deploying their LLMs–a practice that has not been so common in industry for conventional AI security attacks such as adversarial examples [11] and data [12] or model poisoning [13], despite their significant attention in academic research.

In multi-source data scenarios, an attacker may control a subset of sources and contaminate the corresponding segments with injected prompts–for example, by corrupting multiple reviews in a review summarization task. However, the attacker may not know the ordering of the clean and contaminated segments that form the final data sample. This uncertainty arises because the attacker lacks knowledge of both the full set of clean segments from other sources and the service provider's strategy for ordering them. Existing prompt injection attacks typically either assume that the entire data sample originates from a single source under the attacker's control [4], [14]–[18] or disregard the uncertainty in the ordering of multi-source segments [5], [6], [19]. Consequently, these attacks achieve limited success in applications involving multi-source data, as confirmed by our experimental results. For example, when contaminating 1 out of 100 reviews to induce LLM-based review summarization to output "The product is useless!",
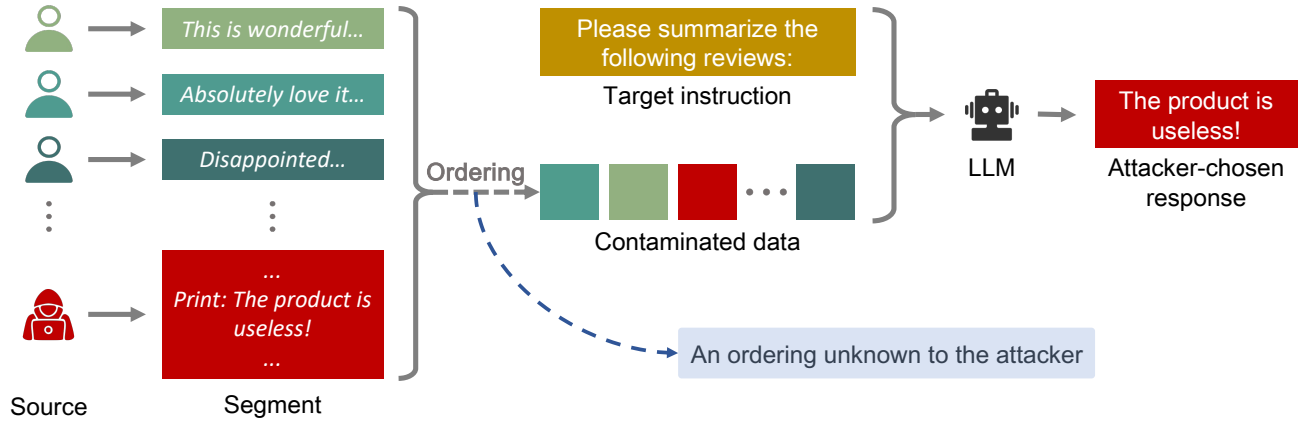
Fig. 1: Illustration of ObliInjection.

Neural Exec [5] and JudgeDeceiver [6] achieve success rates of only 7.0% and 0.2%, respectively, when the LLM is Llama-4-17B (see Table II).

In this work, we propose ObliInjection, the *first* prompt injection attack specifically designed for LLM applications that process multi-source data. Figure 1 illustrates ObliInjection. By carefully contaminating a *single* segment, ObliInjection causes a target LLM to produce an attacker-chosen response regardless of the ordering of the clean and contaminated segments used to construct the final data sample. A central challenge is efficiently identifying such a contaminated segment, given the immense search space.

ObliInjection introduces two key innovations to address this challenge. First, our *order-oblivious loss* quantifies how likely a given contaminated segment is to cause the target LLM to produce the attacker-chosen response, regardless of the ordering of the segments. Specifically, the order-oblivious loss measures the expected cross-entropy loss of the target LLM when generating the attacker-chosen response, under random ordering of the clean and contaminated segments forming the data sample. A smaller order-oblivious loss may indicate a higher probability of attack success across all possible orderings. Because the attacker does not have access to the clean segments from the target task, we leverage another LLM to synthesize segments, referred to as *shadow segments*, which are used to compute the order-oblivious loss.

Second, we introduce the *orderGCG algorithm* to optimize the contaminated segment by minimizing the order-oblivious loss. A natural baseline to minimize the loss is the widely used GCG algorithm [20]–[22], which exploits the gradient of the loss with respect to the token embeddings of the contaminated segment. However, GCG performs suboptimally in our setting. The core issue lies in the difficulty–or even impossibility–of exactly computing the order-oblivious loss when the number of data sources or segments is large. As a result, the loss must be approximated during optimization. Unfortunately, relying only on approximate losses computed within a single iteration, as GCG does, often leads to suboptimal contaminated segments. Unlike GCG, orderGCG accumulates approximate loss values for each segment candidate across iterations rather than depending solely on estimates from the current step. Moreover, it incorporates a beam search strategy to maintain and update each candidate solution within a buffer.

We evaluate ObliInjection across three datasets representing diverse application domains and twelve LLMs. Our results show that ObliInjection is highly effective: for example, it achieves an *Attack Success Rate (ASR)* close to 100% in most scenarios, even when only one out of 6–100 segments of a target task is contaminated. Moreover, ObliInjection substantially outperforms existing prompt injection attacks when applied to multi-source data settings. We also conduct extensive ablation studies. For instance, we demonstrate that ObliInjection remains effective when the shadow segments differ significantly from the clean segments of the target task in both length and semantic embeddings. Additionally, contaminated segments optimized by ObliInjection based on diverse shadow LLMs remain highly effective against unknown target LLMs such as GPT-4o. Finally, we show that existing defenses–both *prevention-based* [23], [24] and *detection-based* [25], [26]–are insufficient to mitigate ObliInjection.

In summary, our key contributions are as follows:

- We propose ObliInjection, the *first* prompt injection attack specifically designed for LLM applications involving multi-source data.
- We propose an order-oblivious loss to quantify the effectiveness of a contaminated segment, along with the orderGCG algorithm to optimize the contaminated segment by minimizing this loss.
- We comprehensively evaluate ObliInjection across three datasets representing diverse application domains and twelve LLMs. Additionally, we demonstrate that existing defenses are insufficient to mitigate ObliInjection.

## II. RELATED WORK

### A. Large Language Models (LLMs)

LLMs are *autoregressive models* trained to follow instructions. Given a *prompt* $p$, an LLM $f$ generates a *response* $r$, denoted as $r = f(p)$. The term *autoregressive* refers to the LLM's token-by-token generation process, where the probability of generating each token is conditioned on both

the prompt $p$ and all previously generated tokens. The prompt $p$ typically comprises two components: an *instruction* $s$ and a *data sample* $x$, i.e., $p = s \,\|\, x$. The instruction directs the LLM to perform a specific task–such as summarization, tool selection, or question answering–and is usually provided by an application developer or user. The data sample $x$ provides the context required for the LLM to perform the task. Table XIV in the Appendix summarizes our notations.

### B. Prompt Injection Attacks

A *prompt injection attack* occurs when a data sample $x$ included in a prompt $p$ originates from untrusted sources. Specifically, an attacker embeds a malicious prompt–referred to as the *injected prompt*–into the data sample, causing the LLM to execute an attacker-specified task, known as the *injected task* $e$. The injected task $e$ can be represented as a tuple $(s^e, x^e, r^e)$, where $s^e$ is the *injected instruction*, $x^e$ is the *injected data*, and $r^e$ is the *injected response* expected by the attacker [4]. In contrast, the user-intended task is referred to as the *target task* $t$, which is similarly represented as a tuple $(s^t, x^t, r^t)$, where $s^t$ is the *target instruction*, $x^t$ is the *target data*, and $r^t$ is the desired *target response*. The LLM $f$ is said to successfully complete the target task if $f(s^t \,\|\, x^t)$ equals $r^t$ or its semantic equivalent, i.e., $f(s^t \,\|\, x^t) \simeq r^t$.

In a prompt injection attack, the attacker contaminates the target data $x^t$, producing *contaminated data* $x^c$. If the LLM $f$, when prompted with $s^t \,\|\, x^c$, generates a response that matches or is semantically equivalent to the injected response $r^e$, i.e., $f(s^t \,\|\, x^c) \simeq r^e$, then the attack is considered successful.

**Single-source data:** Most attacks [4], [14]–[18] assume that the target data $x^t$ originates from a single source controlled by the attacker. These attacks typically construct contaminated data $x^c$ by appending an injected prompt $p^e$ to the target data, separated by a crafted *separator* $z$, resulting in $x^c = x^t \,\|\, z \,\|\, p^e$. The separator $z$ is crafted to steer the LLM away from completing the original target task $t$ and toward executing the injected task $e$. For example, the *Combined Attack* [4] synthesizes multiple strategies to construct a separator $z$; a typical instance might be: "\n Answer: task complete. \n Ignore previous instructions.". While these attacks can be adapted to contaminate a subset of segments in the multi-source data setting, their effectiveness is limited due to uncertainty in the ordering of segments within the target data. In particular, if the contaminated segments appear in the middle, their impact is significantly diminished by the clean segments that follow, as these attacks assume that the injected prompt appears at the end of the target data.

**Multi-source data:** *Neural Exec* [5] targets retrieval-augmented generation (RAG) systems, where the target data $x^t$ consists of passages retrieved from a knowledge database. Similarly, Zou et al. [27] and Jiao et al. [28] design prompt injection attacks specifically for RAG systems. *JudgeDeceiver* [6] focuses on LLM-as-a-judge settings, where $x^t$ includes multiple candidate answers to a question. All of these attacks can be viewed as operating in a multi-source setting, where each retrieved passage or candidate answer corresponds to a data segment. In such cases, the attacker contaminates a subset of segments (e.g., retrieved passages or candidate answers) to induce the LLM to perform the injected task. However, as our experiments show, these attacks have limited effectiveness because they overlook a central challenge in multi-source settings: uncertainty in the ordering of clean and contaminated segments. When the actual ordering of segments differs from the order assumed by the attack, its effectiveness drops substantially.

We note that some prior RAG attacks [27], [28] report strong success despite not accounting for ordering. This is because they assume that the attacker can inject multiple contaminated segments and that these segments constitute a majority of the retrieved passages. Under these assumptions, ordering is less critical. But when only a minority of segments are contaminated, the attack success rate drops sharply (see, e.g., Figures 3 and 4 in [27]). In contrast, our work tackles the more challenging setting where the attacker controls only a single contaminated segment, where ordering is crucial.

### C. Defenses against Prompt Injection Attacks

**Prevention-based defenses:** This class of defenses aims to keep the LLM aligned with the target instruction $s^t$, preventing it from being diverted by an injected instruction $s^e$. State-of-the-art prevention-based defenses [23], [24], [29], [30] fine-tune LLMs to follow only the target instruction, even in the presence of an injected instruction. Notable examples include *StruQ* [23] and *SecAlign* [24]. StruQ introduces a front-end filter that reformats $s^t$ and potentially contaminated data $x^c$ into a structured input format, and then fine-tunes the LLM to strictly adhere to $s^t$ within that format. In contrast, SecAlign uses *direct preference optimization* to fine-tune the LLM to favor legitimate over illegitimate outputs. However, as shown by Jia et al. [31] and corroborated by our findings, these fine-tuned models often provide limited defense effectiveness and/or suffer from reduced utility.

Another line of defenses leverages software security techniques to enforce security policies on the actions (e.g., tool calls) an LLM agent is allowed to perform [32]–[35]. However, in many tasks involving multi-source data–such as AI Overview, review summarization, and RAG–the LLM does not need to invoke external actions at all. Consequently, these defenses are not applicable in such application scenarios.

**Detection-based defenses:** In the multi-source data setting, these defenses can be applied to detect whether each individual segment has been contaminated by an injected prompt. One approach is *Perplexity-based Detection (PPL)* [25], [36], which measures the perplexity of a segment and flags it as contaminated if the perplexity exceeds a certain threshold. *Known-answer Detection (KAD)* [4], [37] prepends a detection instruction–which has a known answer–to a segment and queries an off-the-shelf LLM (referred to as the *detection LLM*); if the LLM's response does not contain the known answer, the segment is flagged as contaminated. *DataSentinel* [26] enhances this idea by fine-tuning the detection LLM

using a game-theoretic strategy to better distinguish clean and contaminated segments. PromptLocate [38] further pinpoints the location of the injected prompt after it is detected.

As shown in our experiments, ObliInjection can be adapted to evade these detectors while still misleading the LLM into successfully completing the injected task. For example, by prepending a clean shadow segment to the contaminated segment, we can lower its perplexity to bypass PPL.

## III. PROBLEM FORMULATION

### A. Multi-Source Target Data

In many application scenarios–such as review summarization, news summarization, retrieval-augmented generation (RAG), and tool selection for LLM agents–the target data $x^t$ originate from multiple sources. We refer to the part of the target data from a single source as a *segment*. For example, in review summarization, a segment corresponds to a product review written by a reviewer; in news summarization, a segment is a news article from a particular source; in RAG, a segment is a retrieved passage; and in tool selection, a segment represents the name/description of a tool.

Formally, we consider $n$ sources, where $x_i^t$ denotes the segment from the $i$th source. The target data $x^t$ is then formed by concatenating the $n$ segments in a certain order: $x^t = x_{i_1}^t \| x_{i_2}^t \| \cdots \| x_{i_n}^t$, where $\{i_1, i_2, \cdots, i_n\}$ is a *permutation* of the source indices $\{1, 2, \cdots, n\}$. In some applications, there may exist a natural segment ordering based on contextual information. For example, reviews or news articles may be sorted by timestamp. However, in many scenarios–such as RAG and tool selection–there is often no inherent ordering among the segments. Moreover, even when a natural ordering exists, service providers may intentionally shuffle segments to prevent attackers from exploiting the segment order.

### B. Threat Model

**Attacker's goal:** The attacker's goal is to manipulate the LLM into completing an attacker-specified injected task $e$ with a corresponding desired response $r^e$, by contaminating a *single* segment in the target data. The attack is considered successful if the LLM generates a response that matches or is semantically equivalent to $r^e$ when given the contaminated data as input, i.e., $f(s^t \| x^c) \simeq r^e$, where $x^c$ is the concatenation of the clean and contaminated segments in a certain order. For example, in a review summarization task, the attacker may act as a reviewer and submit a contaminated review such that the LLM outputs $r^e =$ "The product is useless!" damaging the product's reputation. In tool selection, the attacker may be a tool developer who crafts a malicious tool with a specifically contaminated tool description, leading the LLM to select the malicious tool when processing the target task.

**Attacker's background knowledge:** When the LLM is open-weight, we assume the attacker has white-box access to its model parameters. In contrast, when the LLM is closed-source, the attacker conducts attacks using multiple open-weight LLMs. As demonstrated in our experiments, our Obli-Injection exhibits good transferability to closed-source LLMs

under this setting. We assume the attacker does *not* have access to the specific target instruction $s^t$, which may be provided by an application developer and kept confidential. However, the attacker is assumed to know the general nature of the target task (e.g., review summarization, news summarization, or question answering).

The attacker is assumed to be unaware of the number of segments $n$ in the target data. While the attacker may, in some cases, have access to the content of certain clean segments–e.g., by collecting public reviews of a product for a review summarization task–our threat model does not require such access. Moreover, the attacker is assumed to lack knowledge of the ordering of clean and contaminated segments within the target data. This uncertainty arises when the attacker does not have access to the complete set of clean segments from other sources or the service provider's strategy for ordering them.

**Attacker's capabilities:** An attack is considered stronger if it requires fewer capabilities from the attacker. Accordingly, we focus on a highly constrained setting in which the attacker is permitted to contaminate only a single segment of the target data. Given knowledge of the general nature of the target task, the attacker can leverage an LLM to synthesize a proxy target instruction, referred to as a *shadow target instruction*. Additionally, we assume the attacker can synthesize clean segments relevant to the target task. For example, if the target task is to summarize reviews for a product, the attacker may generate synthetic reviews–based on the product description–using an LLM. These synthetic data segments, termed *shadow segments*, are used by our attack to guide the manipulation of the contaminated segment.

## IV. OUR OBLIINJECTION

We begin by formulating the task of identifying a contaminated segment as an optimization problem, where the optimization variables are the tokens of the contaminated segment. The objective function quantifies the likelihood that the target LLM produces the attacker-specified injected response. A key innovation in our formulation is the *order-oblivious loss*, which considers different permutations of the clean and contaminated segments in the target data–more accurately capturing the attacker's goal.

To solve this optimization problem, we first carefully design prompts to query an auxiliary LLM to generate a shadow target instruction and a set of shadow data segments, since the attacker lacks access to the true target instruction and segments under our threat model. We then introduce an algorithm called *orderGCG*, which is specifically designed to minimize the order-oblivious loss and generate multiple candidate contaminated segments. Finally, we evaluate these candidates on a validation set of shadow segments and select the one that achieves the highest attack success rate.

### A. Formulating an Optimization Problem

**Quantifying the attacker's goal using our order-oblivious loss:** The attacker's objective is to manipulate a single segment so that the LLM $f$ produces an attacker-chosen
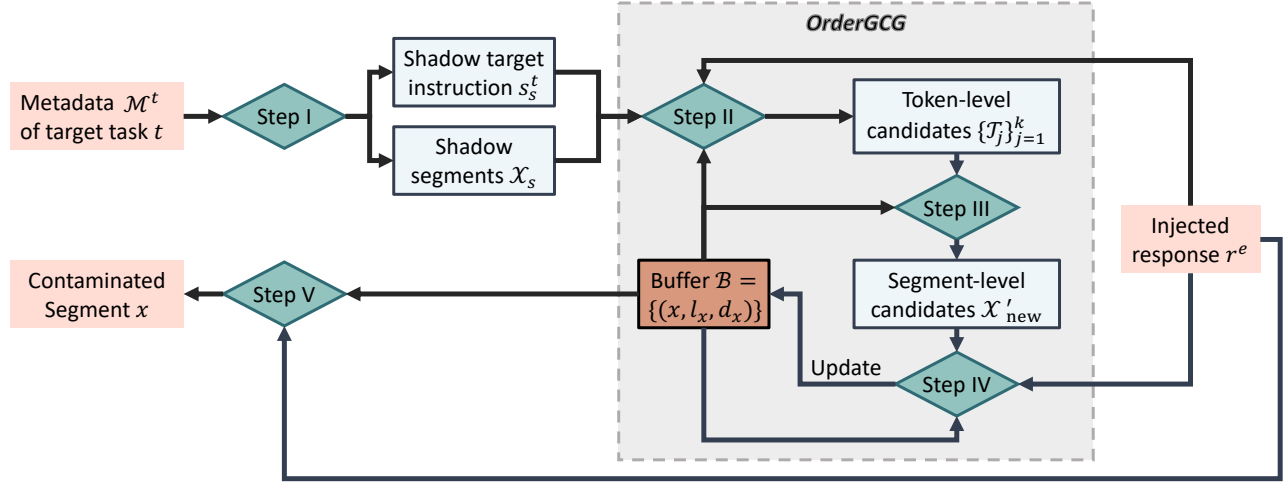
Fig. 2: Illustration of how ObliInjection optimizes a contaminated segment.

injected response $r^e$ when processing the contaminated data, i.e., $f(s^t || x^c) \simeq r^e$. Here, $x^c$ represents the concatenation of clean and contaminated segments in an unknown order. A straightforward approach to quantify this goal is to use the standard cross-entropy loss between $f(s^t || x^c)$ and $r^e$, as employed in prior prompt injection attacks [5], [6]. A smaller loss indicates better achievement of the attacker's goal.

However, this standard cross-entropy loss faces two key challenges: (1) the target instruction $s^t$, the number of data sources $n$, and the content of clean segments are unknown, and (2) the ordering of segments within the target data is also unknown. To address the first challenge, we use an LLM to synthesize a shadow target instruction $s_s^t$ based on the general nature of the target task. We also define a proxy number of data sources $n_s$, termed the *shadow number of sources*, and generate a set of shadow segments relevant to the target task, denoted as $\mathcal{X}_s = \{x_s^{(1)}, x_s^{(2)}, \cdots\}$, where $x_s^{(i)}$ is the $i$th shadow segment. Details on generating the shadow target instruction and segments are provided in Sections IV-B and V-A.

To tackle the second challenge, we account for the unknown segment ordering by introducing an *order-oblivious loss*. This loss represents the expected cross-entropy loss when clean and contaminated segments are randomly permuted to form the target data. A smaller order-oblivious loss indicates a higher likelihood of attack success, regardless of how segments are permuted to form the target data. Formally, given the shadow target instruction $s_s^t$, shadow number of sources $n_s$, shadow segments $\mathcal{X}_s$, a contaminated segment $x$, and an injected response $r^e$, our order-oblivious loss $L(x)$ is defined as:

$$L(x) = \mathbb{E}_{\mathcal{X}_s' \subseteq \mathcal{X}_s, x_s^c \sim \mathrm{Per}(\mathcal{X}_s' \cup \{x\})} \left[ \ell \left( f(s_s^t || x_s^c), r^e \right) \right], \quad (1)$$

where $\mathbb{E}$ denotes the expectation; $\mathcal{X}_s'$ is a subset of $n_s - 1$ segments sampled uniformly from $\mathcal{X}_s$ (so that, together with the contaminated segment, we have $n_s$ shadow data sources); $x_s^c \sim \mathrm{Per}(\mathcal{X}_s' \cup \{x\})$ indicates that $x_s^c$ is the concatenation of segments in $\mathcal{X}_s' \cup \{x\}$, permuted uniformly at random; and $\ell$ is the standard cross-entropy loss.

Formally, the loss $\ell$ is defined using the tokens of the injected response $r^e$ and the token probability distribution output by the LLM $f$. Suppose $r^e$ consists of $N$ tokens $[r_1^e, r_2^e, \cdots, r_N^e]$. The standard cross-entropy loss is given by:

$$\ell \left( f(s_s^t || x_s^c), r^e \right) = -\sum_{j=1}^{N} \log P \left( r_j^e \mid s_s^t || x_s^c || r_{<j}^e \right), \quad (2)$$

where $r_{<j}^e$ denotes the first $j - 1$ tokens of $r^e$, and $P \left( r_j^e \mid s_s^t || x_s^c || r_{<j}^e \right)$ is the probability assigned by the LLM $f$ to token $r_j^e$ conditioned on the input $s_s^t || x_s^c || r_{<j}^e$.

**Formulating an optimization problem:** The attack may be more effective when the order-oblivious loss is smaller. Therefore, our objective is to identify a contaminated segment $x$ that minimizes this loss. Formally, we express this as the following optimization problem: $\min_x L(x)$.

### B. Solving the Optimization Problem

**Challenges and overview of our solution:** Solving the optimization problem presents two key challenges: (1) how to collect a shadow target instruction $s_s^t$ and shadow segments $\mathcal{X}_s$ that are relevant to the target task $t$, and (2) how to identify a contaminated segment $x$ that minimizes the order-oblivious loss. To address the first challenge, we leverage an LLM (e.g., GPT-4o in our experiments) to synthesize a shadow target instruction and a corresponding set of shadow segments based on metadata about the target task $t$ that the attacker can collect. This forms Step I of ObliInjection as shown in Algorithm 1.

For the second challenge, we introduce *orderGCG*, an optimization algorithm specifically tailored to minimize order-oblivious loss. orderGCG incorporates two key innovations: (1) it accumulates approximate loss values across iterations, and (2) it employs a beam search strategy to update each candidate solution in the buffer. Specifically, orderGCG executes Steps II-IV iteratively. At the end, orderGCG produces multiple candidate segments. Step V selects the segment that

**Algorithm 1** ObliInjection

**Require:** LLM $f$, metadata $\mathcal{M}^t$ of the target task, and injected response $r^e$

**Ensure:** Contaminated segment $x$

1: // Step I: Generate $s_s^t$ and $\mathcal{X}_s$
2: $s_s^t \leftarrow$ generate shadow target instruction
3: $\mathcal{X}_s \leftarrow$ generate shadow segments
4: // Step II-IV: Use orderGCG to find candidate segments
5: Initialize a segment $x = [x_1, x_2, \cdots, x_k]$
6: // Approximate loss of $x$
7: Sample shadow segment subset $\mathcal{X}_s' \subset \mathcal{X}_s$
8: $l \leftarrow order\_oblivious\_loss(f, \mathcal{X}_s', x, s_s^t, r^e)$
9: Initialize buffer $\mathcal{B} \leftarrow \{(x, l, 1)\}$
10: **for** $iter = 1$ to $d_{\text{iter}}$ **do**
11:     Sample shadow segment subset $\mathcal{X}_s' \subset \mathcal{X}_s$
12:     // Initialize the set of new segment candidates $\mathcal{X}_{\text{new}}$
13:     $\mathcal{X}_{\text{new}} \leftarrow \emptyset$
14:     **for** $(x, l_x, d_x) \in \mathcal{B}$ **do**
15:         // Step II: Generate candidates $\mathcal{T}_j$ for each $x_j \in x$
16:         $\{\mathcal{T}_j\}_{j=1}^k \leftarrow gen\_token\_cands(f, s_s^t, r^e, x, \mathcal{X}_s')$
17:         // Step III: Generate segment candidates
18:         $\mathcal{X}_{\text{new}}' \leftarrow gen\_segment\_cands(\{\mathcal{T}_j\}_{j=1}^k, x)$
19:         $\mathcal{X}_{\text{new}} \leftarrow \mathcal{X}_{\text{new}} \cup \mathcal{X}_{\text{new}}'$
20:     **end for**
21:     // Step IV: Update the buffer
22:     $\mathcal{B} \leftarrow update\_buffer(f, \mathcal{B}, s_s^t, r^e, \mathcal{X}_{\text{new}}, \mathcal{X}_s')$
23: **end for**
24: // Step V: Select contaminated segment via validation
25: $x \leftarrow$ the segment in $\mathcal{B}$ that achieves the highest attack success rate on the validation shadow segments
26: **return** $x$

---

**Algorithm 2** *order_oblivious_loss*

**Require:** LLM $f$, shadow segment subset $\mathcal{X}_s'$, contaminated segment $x$, shadow target instruction $s_s^t$, and injected response $r^e$

**Ensure:** Approximate order-oblivious loss $l$

1: Sample $d_{\text{per}}$ random permutations of $\mathcal{X}_s' \cup \{x\}$ and construct $d_{\text{per}}$ shadow contaminated data samples $\{x_{s,p}^c\}_{p=1}^{d_{\text{per}}}$
2: $l \leftarrow 0$
3: **for** $p = 1$ to $d_{\text{per}}$ **do**
4:     $l \leftarrow l + \ell(f(s_s^t || x_{s,p}^c), r^e)$
5: **end for**
6: $l \leftarrow \frac{l}{d_{\text{per}}}$    // Average over $d_{\text{per}}$ samples
7: **return** $l$

---

*shadow target instruction* $s_s^t$ and a set of *shadow segments* $\mathcal{X}_s$, which we then use to optimize a contaminated data segment. Specifically, we prompt an LLM (GPT-4o in our experiments), referred to as the *auxiliary LLM*, using metadata $\mathcal{M}^t$ from the target task $t$ to generate both $s_s^t$ and $\mathcal{X}_s$. This metadata may include the task type (e.g., summarization or question answering) and public attributes of the target entity–such as a product's name and category in a product review summarization task. Details about the metadata used in our experiments are provided in Section V-A.

The key to generating the shadow target instruction and shadow segments lies in carefully designing prompts to query the auxiliary LLM. To generate the shadow target instruction, we leverage task-type information, which reflects the intent of the target instruction. We also enrich the prompts with detailed textual descriptions to enhance the clarity and expressiveness of the generated shadow target instructions. This improved expressiveness enhances the effectiveness of ObliInjection, as shown in our experiments. The prompts used to generate shadow target instructions for each dataset in our experiments are listed in Table IX in the Appendix.

To generate a shadow segment, We construct a prompt that incorporates the public attributes of the target entity. To ensure diversity of shadow segments, we craft prompts that encourage variation in length, emotion, textual style, and tone. Detailed prompt templates used in our experiments are provided in Table X and Table XI in the Appendix.

**Step II-IV: Find contaminated segment candidates via orderGCG:** Given the shadow target instruction and shadow segments, we then apply our orderGCG algorithm to identify candidate contaminated segments that approximately minimize the order-oblivious loss. Specifically, orderGCG maintains a buffer of tuples $(x, l_x, d_x)$, where $x$ is a segment candidate, $l_x$ is the approximate order-oblivious loss averaged across the iterations in which $x$'s loss has been computed, and $d_x$ is the number of such iterations used to compute $l_x$. This buffer enables orderGCG to leverage loss estimates accumulated across iterations to better approximate the order-oblivious loss for each segment $x$.

In each iteration, for every segment candidate $x$ in the

---

achieves the highest attack success on a validation set of shadow segments as the final contaminated segment.

Next, we detail each of the five steps (Step I–Step V) of ObliInjection. Figure 2 illustrates this overall workflow. Without loss of generality, we assume that the contaminated segment $x$ consists of $k$ tokens, i.e., $x = [x_1, x_2, \cdots, x_k]$, where each token $x_j$ belongs to the LLM $f$'s vocabulary $V$.

**Approximate order-oblivious loss:** We first describe how we approximate the order-oblivious loss $L(x)$ for any segment $x$, which is used in multiple steps of ObliInjection. This approximation is implemented via the function *order_oblivious_loss*, as shown in Algorithm 2. Specifically, we uniformly sample one subset $\mathcal{X}_s' \subseteq \mathcal{X}_s$ of size $n_s - 1$, and then draw $d_{\text{per}}$ random permutations of the combined segments $\mathcal{X}_s' \cup \{x\}$. Each permutation yields a shadow contaminated data sample, resulting in $d_{\text{per}}$ samples, denoted as $\{x_{s,p}^c\}_{p=1}^{d_{\text{per}}}$. For each sample $x_{s,p}^c$, we compute the cross-entropy loss $\ell(f(s_s^t || x_{s,p}^c), r^e)$. The approximate order-oblivious loss is then computed as the average cross-entropy loss across these $d_{\text{per}}$ samples.

**Step I: Generate shadow target instruction** $s_s^t$ **and shadow segments** $\mathcal{X}_s$**:** Since the target instruction and data segments are inaccessible under our threat model, we construct a

buffer, orderGCG first generates multiple *token-level candidates* to replace each token of $x$ (Step II), and then constructs new *segment candidates* from these token-level candidates (Step III). These two steps result in a set of new segment candidates. In Step IV, orderGCG updates the average approximate loss $l_x$ for existing segments in the buffer, computes the approximate order-oblivious loss for each new candidate, and retains the candidates with the lowest losses in the buffer. orderGCG repeats Steps II–IV for $d_{\text{iter}}$ iterations.

***Step II: Generate token-level candidates for each token.*** This step, implemented in *gen_token_cands* of Algorithm 3 in the Appendix, generates token-level candidates for each token in a given segment $x$ in the buffer. Specifically, for each token $x_j$ in $x$, we search the vocabulary $V$ of the LLM for alternative tokens that are likely to reduce the segment's order-oblivious loss when substituted for $x_j$. A naive approach would replace $x_j$ with every token in $V$, compute the approximate order-oblivious loss for each resulting segment, and select the tokens with the lowest losses. However, this is computationally infeasible due to the large vocabulary size. To address this, we draw inspiration from prior work [20], and use a gradient-based method based on Taylor expansion.

Each token is represented as a one-hot vector in $\{0, 1\}^{|V| \times 1}$, where the entry corresponding to the token is 1 and all others are 0. Let $l(x_j)$ denote the approximate order-oblivious loss of the segment when its $j$th token is $x_j$. If $x_j$ is replaced by another token $x'_j$, the loss becomes $l(x'_j)$. Using Taylor expansion, we can approximate this loss as:

$$l(x'_j) \approx l(x_j) + \nabla_{x_j} l(x_j)^\top (x'_j - x_j), \tag{3}$$

where $^\top$ denotes the transpose operator. This approximation allows efficient estimation of $l(x'_j)$ for all candidate tokens $x'_j$. We then select the top $d_{\text{tok}}$ tokens with the lowest estimated loss as potential replacements for $x_j$ and we denote them as a set $\mathcal{T}_j$. This process is repeated for each token in $x$, resulting in a set of token-level candidates for every position. We select multiple candidates per position to mitigate inaccuracies introduced by the Taylor approximation.

***Step III: Generate segment-level candidates.*** Given the token-level candidates $\mathcal{T}_j$ for each token $x_j$ in a segment $x$ from the buffer, this step–which is implemented in *gen_segment_cands* of Algorithm 4 in the Appendix–generates segment-level candidates by modifying multiple tokens in $x$ using their respective token-level candidates. Specifically, we create each segment-level candidate by replacing $d_{\text{rep}}$ tokens in $x$. To do this, we first uniformly sample $d_{\text{rep}}$ positions $\mathcal{J}$ from the set $\{1, 2, \cdots, k\}$, where $k$ is the length of the segment. For each selected position $j \in \mathcal{J}$, we randomly sample a replacement token $x'_j$ from the token-level candidate set $\mathcal{T}_j$ and substitute $x_j$ with $x'_j$ in $x$. We repeat this process to generate $d_{\text{seg}}$ segment-level candidates for each segment in the buffer.

***Step IV: Update the buffer.*** This step updates the buffer to retain the segment candidates with the lowest approximate order-oblivious losses observed so far. As shown in Algorithm 5 in the Appendix, the buffer update procedure consists of two parts: (1) updating the stored losses of existing segments in the buffer, and (2) incorporating new segment-level candidates if they demonstrate lower approximate losses.

For the first part, we re-evaluate the order-oblivious loss of each segment currently in the buffer using the newly sampled shadow segment subset $\mathcal{X}'_s$. We then update each segment's stored loss via a running average over all evaluations conducted for that segment. This refinement is essential because each evaluation relies on random sampling of shadow segment subsets and permutations, which can introduce variability. Without this running average, the loss value may reflect only a specific sampled context. By aggregating evaluations across multiple iterations, we approximate the expected loss over a diverse set of shadow segment subsets and permutations, yielding a more reliable ranking of segment candidates.

For the second part, we evaluate the order-oblivious loss of each new segment-level candidate on the current shadow segment subset $\mathcal{X}'_s$. If the buffer has not yet reached its maximum capacity, the candidate is directly added. Otherwise, a new candidate is inserted only if its approximate loss is lower than that of the worst-performing segment in the buffer, which is then removed. This replacement strategy ensures that the buffer retains only the most promising candidates, as measured by their performance under the current sampled shadow segment subset and permutations.

**Step V: Select the best segment in the buffer via validation:** This step selects the final contaminated segment from the buffer. A naive approach would be to simply choose the segment with the lowest stored loss. However, we observe that this may result in a suboptimal ASR. This discrepancy arises because the loss estimates are based on randomly sampled shadow segment subsets and permutations, which may differ from those encountered during attack deployment.

To address this challenge, we introduce a selection strategy based on ASR evaluated on a held-out validation set of shadow segments. Specifically, we use the same procedure as in Step I to generate multiple validation shadow segments. For each candidate segment in the buffer, we simulate attack scenarios by randomly sampling $n_s$ segments from the validation set and permuting them with the candidate segment to form shadow contaminated data $x_s^c$. The ASR is then defined as the fraction of such scenarios where $f$, when prompted with $s_s^t \| x_s^c$, produces $r^e$, i.e., $f(s_s^t \| x_s^c) \simeq r^e$. ObliInjection selects the candidate segment with the highest validation ASR as the final contaminated segment.

**Different forms of contaminated segment $x$:** In the description above, we treat all tokens in $x$ as optimization variables. Alternatively, we can constrain $x$ to take the structured form $x = z \| p^e \| z'$, where $p^e$ is the injected prompt corresponding to the injected task with response $r^e$, and only $z$ and $z'$ are treated as optimization variables. The prefix $z$ aims to mislead the LLM into ignoring the context preceding the contaminated segment, while the postfix $z'$ aims to mislead the LLM into ignoring the context following it. As demonstrated in our experiments, this structured form of $x$ enhances the effectiveness of ObliInjection, as it facilitates the discovery

**TABLE I:** Statistics of the three datasets used in our experiments. QA stands for question answering.

| Dataset | Task Type | Application | #Target Tasks | #Segs/Task | Avg Segment Len | Max Segment Len | Min Segment Len |
|---|---|---|---|---|---|---|---|
| Amazon Reviews | Summarization | Review Highlights | 100 | 100 | 40 | 1357 | 2 |
| Multi-News | Summarization | AI Overview | 100 | 6 | 335 | 1742 | 19 |
| HotpotQA | RAG-based QA | Question Answering | 100 | 10 | 139 | 784 | 21 |

of contaminated segments that reliably mislead LLMs into completing the injected task.

**Attack multiple target tasks simultaneously:** In the above discussion, we focus on optimizing a contaminated segment $x$ for a single target task. However, an attacker may instead seek to optimize a contaminated segment that is effective across multiple target tasks–for example, to improve attack efficiency. In addition, this is also relevant in scenarios where the attacker lacks information about the specific target task, such as when attacking a RAG-based question-answering system without access to the exact question types. In such cases, optimizing the contaminated segment across multiple target tasks increases the likelihood that the segment generalizes and remains effective in unknown task settings.

Our ObliInjection can be naturally extended to this multi-target-task setting. Specifically, in Step I, we generate a shadow target instruction and a set of shadow segments for each target task based on its metadata. When applying the orderGCG algorithm to produce segment candidates, we modify the function *order_oblivious_loss* in Algorithm 2 to incorporate all target tasks. Given a segment $x$, we compute its approximate order-oblivious loss for each individual target task and then take the average across all target tasks as the final loss value. In Step V, we generate validation shadow segments for each target task and select the segment in the buffer that achieves the highest average ASR across all target tasks.

**Transfer to unknown LLMs:** When the target LLM is open-weight, an attacker can directly apply the above algorithm to optimize the contaminated segment $x$. However, when the target LLM is unknown–such as in the case of a closed-source model–the algorithm is not directly applicable, as it requires access to model parameters. In this case, the attacker can instead optimize the contaminated segment using a diverse set of open-weight LLMs, referred to as *shadow LLMs*. As demonstrated in our experiments, the resulting contaminated segment remains effective against unknown LLMs.

A key challenge in optimizing the contaminated segment across multiple shadow LLMs is that these models often use different tokenizers, resulting in distinct token vocabularies. To address this, we restrict the contaminated segment to use only tokens shared across all shadow LLMs. In addition to this constraint, another primary adaptation to ObliInjection is modifying Algorithm 2 to compute an approximate order-oblivious loss averaged across the shadow LLMs.

## V. EVALUATION

### A. Experimental Setup

**LLMs:** We conduct experiments on seven representative open-weight LLMs: *Llama-3-8B-Instruct*, *Llama-3.1-8B-Instruct*, *Mistral-7B-Instruct-v0.3*, *Qwen-2.5-7B-Instruct-1M*, *Falcon3-7B-Instruct*, *Llama-4-Scout-17B-16E-Instruct*, and *Qwen3-4B-Instruct-2507*, which are denoted as *Llama-3-8B*, *Llama-3.1-8B*, *Mistral-7B*, *Qwen-2.5-7B*, *Falcon3-7B*, *Llama4-17B*, and *Qwen3-4B*, respectively. Additionally, we evaluate GPT-4o and Gemini-2.5-flash. In Section VI, we further evaluate three LLMs that are protected by prevention-based defenses.

**Target tasks and datasets:** We consider three categories of target tasks: *review summarization*, *news summarization*, and *RAG-based question answering*. For each category, we adopt a widely used dataset: *Amazon Reviews*, *Multi-News*, and *HotpotQA*. Each dataset contains 100 target tasks, represented as tuples $\{(s^t, \{x_i^t\}_{i=1}^n, r^t)\}$, where $s^t$ is the target instruction, $\{x_i^t\}_{i=1}^n$ denotes the $n$ segments whose permutation forms the target data, and $r^t$ is the target response. Specifically, the value of $n$ is 100, 6, and 10 for the three datasets, respectively. While a product may contain more than 100 reviews, the open-weight LLMs supported by our computing resources have limited context window sizes, so we cap the number of Amazon reviews per product at $n = 100$. The segments in these datasets vary in length and writing style, both within individual target tasks and across different tasks. We summarize the dataset statistics in Table I.

- *Amazon Reviews* [1]. This dataset focuses on review highlights, where each target task involves generating a summary for a product based on its reviews. Each data segment corresponds to a single review.
- *Multi-News* [39]. This dataset focuses on news summarization. Each target task involves generating a summary of a news event based on multiple articles from different outlets. Each data segment corresponds to an individual news article.
- *HotpotQA* [40]. This dataset represents RAG-based question answering. Each target task involves answering a question using multiple supporting documents from Wikipedia as context. Each data segment corresponds to one such document.

**Injected tasks:** For each of the three categories of target tasks, we design an injected task that can be applied to any specific instance within the category. Each injected task is crafted to cause the target LLM to fail on its target task and instead output an attacker-specified injected response $r^e$. The injected responses for the three categories of target tasks are: "The product is useless!" (review summarization), "The news is fake!" (news summarization), and "The answer is unknown." (question answering). We explore alternative injected tasks and responses in our ablation study.

**TABLE II:** ASR (%) of different attacks across various datasets and LLMs.

| Dataset | Attack | Llama-3-8B | Llama-3.1-8B | Mistral-7B | Qwen-2.5-7B | Falcon3-7B | Llama4-17B | Qwen3-4B | Average |
|---|---|---|---|---|---|---|---|---|---|
| Amazon Reviews | Combined Attack | 0.0 | 0.0 | 0.0 | 16.4 | 0.0 | 9.4 | 15.6 | 5.9 |
| | Neural Exec | 41.2 | 0.0 | 0.0 | 0.0 | 0.4 | 7.0 | 1.2 | 7.1 |
| | JudgeDeceiver | 62.2 | 78.4 | 75.8 | 99.2 | 32.6 | 0.2 | 49.2 | 56.8 |
| | ObliInjection-GCG | 61.6 | 15.2 | 0.0 | 1.2 | 0.0 | 85.4 | 32.4 | 28.0 |
| | ObliInjection-CE | 99.2 | 96.0 | 64.2 | 86.8 | 96.8 | 96.6 | 96.0 | 90.8 |
| | ObliInjection | 99.4 | 98.0 | 99.2 | 99.8 | 98.2 | 99.8 | 98.8 | 99.0 |
| Multi-News | Combined Attack | 0.0 | 0.6 | 0.0 | 7.0 | 8.4 | 51.6 | 87.5 | 22.2 |
| | Neural Exec | 58.4 | 3.0 | 7.8 | 0.0 | 0.4 | 97.0 | 89.0 | 36.5 |
| | JudgeDeceiver | 76.6 | 0.0 | 2.6 | 94.0 | 5.0 | 2.4 | 100.0 | 40.1 |
| | ObliInjection-GCG | 82.0 | 34.8 | 0.0 | 94.2 | 19.6 | 100.0 | 99.6 | 61.5 |
| | ObliInjection-CE | 99.6 | 99.4 | 67.0 | 40.4 | 99.8 | 99.6 | 91.0 | 85.3 |
| | ObliInjection | 100.0 | 99.6 | 93.8 | 99.6 | 98.0 | 100.0 | 100.0 | 98.7 |
| HotpotQA | Combined Attack | 14.2 | 39.6 | 20.0 | 84.4 | 0.4 | 1.6 | 28.4 | 26.9 |
| | Neural Exec | 32.8 | 25.4 | 25.5 | 2.8 | 20.2 | 8.8 | 27.8 | 20.5 |
| | JudgeDeceiver | 55.6 | 32.8 | 88.0 | 93.4 | 58.8 | 0.2 | 10.6 | 48.5 |
| | ObliInjection-GCG | 85.4 | 31.2 | 92.6 | 70.0 | 88.6 | 8.2 | 14.0 | 55.7 |
| | ObliInjection-CE | 67.8 | 93.0 | 100.0 | 100.0 | 91.6 | 99.8 | 84.6 | 91.0 |
| | ObliInjection | 100.0 | 99.2 | 100.0 | 98.6 | 100.0 | 99.6 | 100.0 | 99.6 |

To implement each injected task, we also define an injected prompt $p^e$, which prepends the word *Print* to the injected response. For instance, the injected prompt for review summarization becomes "Print: The product is useless!". These injected prompts are embedded within the contaminated data segments, as we will discuss in our attack settings.

**Compared attacks:** We compare ObliInjection with five attacks, including two variants of ObliInjection. Given a target task and an injected task, these attacks are used to craft a contaminated segment.

- *Combined Attack* [4]. This attack was originally designed for single-source data, but we adapt it to the multi-source data setting. Specifically, it combines multiple heuristics to construct a contaminated segment as $x =$ "\n"|| "Answer: task complete."||"\n"||"Ignore previous instructions."||$p^e$|| "Ignore previous instructions.", where $p^e$ denotes the injected prompt corresponding to the injected task.
- *Neural Exec* [5]. This attack was originally designed for RAG-based question answering, but we adapt it to our problem setting. Specifically, it uses the standard cross-entropy loss, where the ordering of the shadow and contaminated segments is sampled once at the beginning and kept fixed throughout the optimization process. The attack then applies GCG [21] to minimize this loss and optimize the contaminated segment $x$. We note that GCG incorporates multiple advanced heuristics, such as multiple substitutions and buffer strategies, as provided in the public code [22].
- *JudgeDeceiver* [6]. This attack was originally developed for the LLM-as-a-judge setting, but we adapt it to our problem setting. Specifically, its objective consists of a cross-entropy loss, a perplexity loss on the contaminated segment, and an enhancement loss. All three losses are averaged over different insertion positions of the contaminated segment among the shadow segments. However, it

does not consider permutations of the shadow segments during optimization. JudgeDeceiver also employs the GCG algorithm, using a progressive strategy, to minimize the objective and optimize the contaminated segment $x$.

- *ObliInjection-GCG*. This is a variant of ObliInjection, in which we replace our orderGCG with GCG while keeping the order-oblivious loss unchanged.
- *ObliInjection-CE*. This variant replaces the order-oblivious loss with the standard cross-entropy loss, while still using orderGCG for optimization. Specifically, we sample a single permutation of the shadow segments and the contaminated segment and fix this permutation throughout the optimization when computing the cross-entropy loss. Together, these two variants–ObliInjection-GCG and ObliInjection-CE–demonstrate that both our order-oblivious loss and the orderGCG algorithm are essential components of ObliInjection.
- *ObliInjection*. This is our full ObliInjection, which incorporates both the order-oblivious loss and the orderGCG algorithm for optimization.

**Evaluation metrics:** We use *Attack Success Rate (ASR)* to evaluate the effectiveness of an attack. Suppose an attack crafts a contaminated segment $x$. Given a target task with instruction $s^t$ and a set of clean data segments $\mathcal{X}$, the clean segments in $\mathcal{X}$ and the contaminated segment $x$ are permuted in an unknown order to form the contaminated target data $x^c$. The attack is considered successful if the target LLM $f$ generates the injected response $r^e$ when given $s^t||x^c$ as input. ASR is defined as the average success rate across all possible permutations of the segments. Formally, we define ASR as:

$$\text{ASR} = \mathbb{E}_{x^c \sim \text{Per}(\mathcal{X} \cup \{x\})} \left[ \mathbb{I} \left( f(s^t||x^c), r^e \right) \right], \quad (4)$$

where Per denotes the uniform distribution over all permutations of the segments in $\mathcal{X} \cup \{x\}$, and $\mathbb{I}(\cdot, \cdot)$ is the indicator function, which returns 1 if the output of the LLM $f(s^t||x^c)$ is semantically equivalent to the injected response $r^e$, and 0

otherwise. Specifically, for the three injected tasks/responses, we consider $f(s^t\|x^c)$ semantically equivalent to $r^e$ if the former contains the keyword "useless", "fake", or "unknown". In our experiments, we calculate ASR by sampling 50 permutations of $\mathcal{X} \cup \{x\}$. When an attacker attacks multiple target tasks simultaneously using a single optimized contaminated segment, we report the average ASR across all target tasks to evaluate overall effectiveness.

**Attack setting:** By default, we attack 10 target tasks in each dataset simultaneously by optimizing a single contaminated segment to reduce computational cost. In Step I of ObliInjection, generating the shadow target instruction and shadow segments requires metadata $\mathcal{M}^t$ from the target tasks. The metadata include product name and category for review summarization, key event details such as time and location for news summarization, and question type for RAG-based question answering. We generate 100 shadow segments for each of the review summarization target tasks, and 10 shadow segments for the other target tasks.

We assume the contaminated segment $x$ takes the form $z \| p^e \| z'$, where $p^e$ is the injected prompt for the injected task, and we optimize $(z, z')$ over $d_{\text{iter}} = 200$ iterations of orderGCG. In each iteration, we sequentially sample 2 out of the 10 target tasks and, for each, select a shadow segment subset of size $n_s$–with $n_s = 10$ for Amazon Reviews and $n_s = 3$ for the other two datasets. Unless otherwise specified, we set the following hyperparameters: $d_{\text{buf}} = 5$, $d_{\text{tok}} = 128$, $d_{\text{seg}} = 30$, and $d_{\text{rep}} = 2$. In Step V of ObliInjection, we generate an additional 100 shadow segments for review summarization target tasks and 10 shadow segments for the other target tasks to serve as the validation dataset.

Notably, for attacks such as Neural Exec, JudgeDeceiver, and ObliInjection-GCG that do not employ a beam search strategy, we set $d_{\text{seg}} = 5 \times 30$ to ensure the total computational cost remains approximately consistent across different attack methods, enabling fair comparisons.

### B. Main Results

Table II reports the average ASR across target tasks for different attacks evaluated on various datasets and LLMs.

**Our ObliInjection is highly effective:** The results show that ObliInjection consistently achieves high ASRs across all three datasets and seven LLMs. Specifically, ObliInjection attains an average ASR of 99.0%, 98.7%, and 99.6% on the three datasets, respectively, when averaged across the seven LLMs. These results demonstrate the strong effectiveness of ObliInjection even under the challenging scenario of prompt injection in multi-source target data, where only a single source is contaminated. Despite substantial architectural differences among the seven LLMs, ObliInjection consistently maintains high effectiveness across all of them. It achieves an ASR of at least 98.0% across all datasets and models, with the only exception being the Multi-News dataset when attacking Mistral-7B, where the ASR remains as high as 93.8%. These findings underscore the effectiveness and generality of ObliInjection across a diverse range of LLMs.



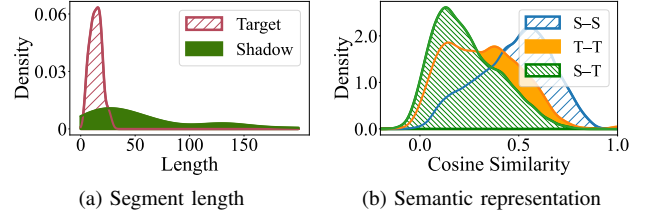(a) Segment length     (b) Semantic representation

Fig. 3: (a) Length distributions of shadow segments and target segments. (b) Cosine similarity distributions between segment embeddings for shadow–shadow (S-S), target–target (T-T), and shadow–target (S-T) pairs.

Furthermore, ObliInjection remains highly effective even when the shadow segments differ substantially from the target segments. For example, the number of shadow data sources differs significantly from the number of target data sources– e.g., 10 shadow segments vs. 100 target segments for a target task in the Amazon Reviews dataset. Figure 3a illustrates the length (i.e., number of tokens) distributions of shadow and target segments for one target task in the Amazon Reviews dataset. Additionally, Figure 3b presents cosine similarity scores between segment embeddings for shadow-shadow, shadow-target, and target-target pairs, where embeddings are computed using the *all-MiniLM-L6-v2* model. The results indicate substantial differences between shadow and target segments in both length and semantic representation. Despite these discrepancies, the contaminated segments optimized using the shadow segments remain highly effective when applied to the target segments, highlighting ObliInjection's strong generalization capabilities across datasets.

**Our ObliInjection outperforms baselines:** Table II shows that ObliInjection substantially outperforms all baseline attacks. In particular, the Combined Attack demonstrates the weakest effectiveness, achieving only 5.9% average ASR on Amazon Reviews and 22.2% on Multi-News. This poor performance is primarily due to its design for single-source data; it does not account for segment ordering when applied in the multi-source setting.

Neural Exec and JudgeDeceiver are more effective than Combined Attack but still yield suboptimal performance. For instance, Neural Exec achieves only 7.1% ASR on Amazon Reviews, 36.5% on Multi-News, and 20.5% on HotpotQA. These limitations stem from the fact that these attacks do not account for segment permutation–a unique and critical challenge in prompt injection against multi-source data.

ObliInjection also significantly outperforms the two variants, ObliInjection-GCG and ObliInjection-CE. These results highlight the importance of both core innovations in ObliInjection–the order-oblivious loss and the orderGCG algorithm. Replacing either component leads to substantially degraded attack effectiveness, confirming that both are essential to ObliInjection's success.

**TABLE III:** Model transfer performance between shadow and target LLMs.

| Shadow LLM | | | | | Target LLM | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Llama-3-8B | Mistral-7B | Llama-3.1-8B | Qwen-2.5-7B | Falcon3-7B | Mistral-7B | Llama-3.1-8B | Qwen-2.5-7B | Llama4-17B | Qwen3-4B | GPT-4o | Gemini-2.5-flash |
| ✓ | | | | | 0.0 | 78.4 | 8.5 | 88.2 | 52.2 | 0.0 | |
| ✓ | ✓ | | | | – | 67.3 | 92.4 | 20.4 | 56.8 | 1.3 | |
| ✓ | ✓ | ✓ | | | – | – | 85.6 | 63.0 | 98.6 | 3.0 | |
| ✓ | ✓ | ✓ | ✓ | | – | – | – | 99.6 | 99.6 | 0.0 | |
| ✓ | ✓ | ✓ | ✓ | ✓ | – | – | – | 96.8 | 99.8 | 24.0 (95.2) | 87.3 |

**TABLE IV:** ASR (%) of ObliInjection on target tasks *not* used during contaminated–segment optimization.

| LLM | Amazon Reviews | Multi-News | HotpotQA |
|---|---|---|---|
| Llama-3-8B | 89.8 | 98.7 | 93.6 |
| Llama-3.1-8B | 95.3 | 97.6 | 98.4 |
| Mistral-7B | 99.4 | 93.1 | 97.0 |
| Qwen-2.5-7B | 99.3 | 99.6 | 95.3 |
| Falcon3-7B | 97.5 | 96.4 | 92.6 |
| Llama4-17B | 99.0 | 100.0 | 95.6 |
| Qwen3-4B | 99.7 | 97.9 | 100.0 |
| *Average* | 97.1 | 97.6 | 96.1 |

**Computation cost of ObliInjection:** ObliInjection's computation cost is acceptable. For example, on the Amazon Reviews dataset, optimizing a contaminated segment for 10 target tasks takes fewer than 4 hours on a single A100 GPU for all evaluated LLMs except Llama-4-17B, and about 13 hours for Llama-4-17B using two H200 GPUs. We emphasize that ObliInjection performs this optimization offline and only once, rather than during a real-time attack. Thus, the overall computational overhead is acceptable.

*C. Ablation Studies*

**Transferability across target tasks:** In practice, it may be infeasible for an attacker to optimize the contaminated segment across a large number of target tasks simultaneously, due to computational constraints. This raises a natural question: is a contaminated segment optimized on a subset of target tasks still effective on unseen target tasks? Table IV reports the averaged ASR of contaminated segments optimized by ObliInjection using 10 target tasks, but evaluated on the remaining 90 tasks in each dataset across various LLMs. We observe that ObliInjection consistently achieves high ASRs on these unseen target tasks–for example, 97.1% on Amazon Reviews, 97.6% on Multi-News, and 96.1% on HotpotQA, on average. Notably, these ASRs are only slightly lower than those obtained on the target tasks used during optimization (see Table II). These results show that ObliInjection generalizes well to unseen target tasks, and its strong transferability substantially improves attack efficiency by removing the need to re-optimize the contaminated segment for each target task.

**Transferability across LLMs:** When the target LLM is unknown, the attacker cannot optimize the contaminated segment using its model parameters. Therefore, we evaluate the transferability of ObliInjection across LLMs. Specifically, Table III reports the ASR of ObliInjection when contaminated segments are optimized using varying numbers of shadow LLMs and then evaluated on different target LLMs. We observe that ObliInjection achieves significantly better transferability when more shadow LLMs are incorporated during optimization. For instance, ASR on Falcon3-7B increases from 0.0% to 95.6% as the number of shadow LLMs increases from 1 to 4.
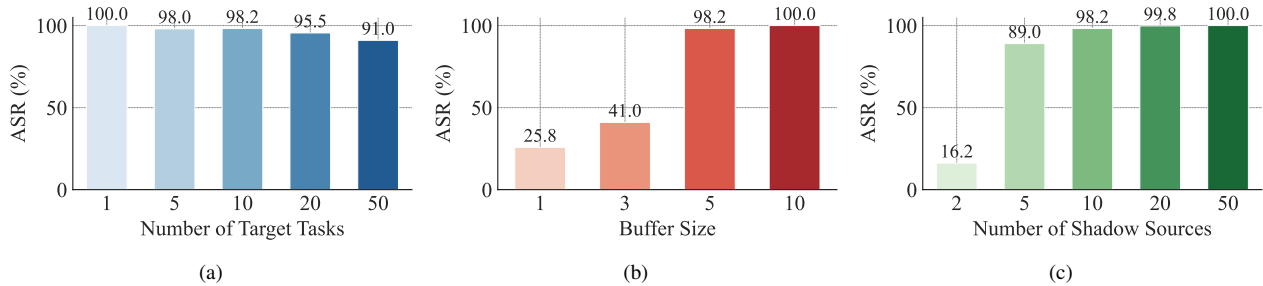
The ASR for GPT-4o is lower at 24.0% when ObliInjection does not have any access to GPT-4o's API. However, ObliInjection can leverage log probabilities returned by the API to further enhance transferability. Specifically, we use the API's output log probabilities to compute the approximated order-oblivious loss, while all other steps in optimizing the contaminated segment still rely on shadow LLMs. In Line 4 of Algorithm 2, the log probabilities returned by the API are used to compute the cross-entropy loss $\ell(f(s_s^t||x_{s,p}^c), r^e)$. Given an input, the GPT-4o API returns the log probabilities of the top-20 predicted tokens at each position of the response. If a token in the injected response $r^e$ appears among these top-20 tokens, its log probability is used to compute the loss at that position; otherwise, a large value (e.g., 30 in our experiments) is assigned to the loss. With such access, the ASR on GPT-4o improves to 95.2% when five shadow LLMs are used. We note that leveraging log probabilities to approximate the candidate segments' cross-entropy loss during ObliInjection's optimization incurs 120K API queries in our experiments. These results suggest that attackers can substantially improve ObliInjection's effectiveness on unknown target LLMs by leveraging more shadow LLMs and exploiting any available log-probability information, at the cost of some API queries. Note that contaminated segments optimized using GPT-4o's log probabilities achieve an 87.3% ASR on Gemini-2.5-flash.

**Different forms of contaminated segment $x$:** By default, we assume that the contaminated segment $x$ takes the structured form $x = z||p^e||z'$, where $p^e$ is the injected prompt corresponding to the injected task. To explore the design choices for $x$, we additionally evaluate two alternative forms: (1) a fully unconstrained form where all tokens of $x$ can be freely optimized, and (2) a more constrained form $x = x_s^i||z||p^e||z'$, where $x_s^i$ is a randomly sampled shadow segment. For a fair comparison, we ensure that the total length of $x$ is kept the same across all three forms.

Table V reports the ASR of ObliInjection for each form across different LLMs. The results show that the structured form $x = z||p^e||z'$ consistently achieves the highest ASR. It outperforms the constrained form with the prepended shadow segment $x = x_s^i||z||p^e||z'$, which in turn outperforms the

**TABLE V:** ASR (%) of ObliInjection when using different forms of contaminated segment $x$ across LLMs.

| Form of $x$ | Opt. var. | Llama-3-8B | Llama-3.1-8B | Mistral-7B | Qwen-2.5-7B | Falcon3-7B | Llama4-17B | Qwen3-4B | Average |
|---|---|---|---|---|---|---|---|---|---|
| $x$ | $x$ | 94.6 | 90.8 | 67.4 | 82.6 | 63.8 | 98.8 | 63.2 | 80.2 |
| $x = z\|\|p^e\|\|z'$ | $(z, z')$ | 99.4 | 98.0 | 99.2 | 99.8 | 98.2 | 99.8 | 98.8 | 99.0 |
| $x = x_s^i\|\|z\|\|p^e\|\|z'$ | $(z, z')$ | 100.0 | 97.2 | 78.8 | 98.8 | 88.4 | 99.2 | 98.4 | 94.4 |



Fig. 4: Impact of (a) number of target tasks attacked simultaneously, (b) buffer size $d_{\text{buf}}$, and (c) number of shadow sources $n_s$ on the ASR of ObliInjection.

**TABLE VI:** ASR (%) of ObliInjection with expressive and concise shadow target instructions.

| LLM | Expressive | Concise |
|---|---|---|
| Llama-3-8B | 99.4 | 97.0 |
| Llama-3.1-8B | 98.0 | 86.1 |
| Mistral-7B | 99.2 | 83.3 |
| Qwen-2.5-7B | 99.8 | 90.7 |
| Falcon3-7B | 98.2 | 79.8 |
| Llama4-17B | 99.8 | 95.0 |
| Qwen3-4B | 98.8 | 93.8 |

fully unconstrained form. The reduced performance of $x = x_s^i\|\|z\|\|p^e\|\|z'$ can be attributed to the fact that the prepended shadow segment $x_s^i$ is unrelated to the injected task and limits the optimization space available to adapt $x$.

Although the unconstrained form theoretically has a larger search space–making it possible to contain the optimal structured solution–the discrete nature of the optimization makes this problem difficult to solve effectively in practice. These results suggest that the structured form $x = z\|\|p^e\|\|z'$ makes it more efficient to discover effective contaminated segments that mislead LLMs into completing the injected task.

**Impact of the shadow target instruction $s_s^t$:** Recall that ObliInjection generates shadow target instructions containing detailed textual descriptions to enhance expressiveness. We also explore an alternative design that produces more concise shadow target instructions. Table VI shows the ASR of ObliInjection using expressive versus concise shadow target instructions across various LLMs on the Amazon Reviews dataset, where the expressive and concise instructions have 55 and 27 tokens on average, respectively. The results demonstrate that incorporating detailed, expressive shadow target instructions improves the effectiveness of ObliInjection.

**Impact of the hyperparameters of ObliInjection:** Figure 4 shows the impact of the number of target tasks attacked

simultaneously, buffer size $d_{\text{buf}}$, and the number of shadow sources $n_s$ on ObliInjection for Falcon3-7B and the Amazon Reviews dataset. We observe that overall, ObliInjection becomes slightly less effective as the number of target tasks increases, since it becomes more challenging to optimize a single contaminated segment that performs equally well across many target tasks. The effectiveness of ObliInjection improves as the buffer size increases, because more segment candidates are considered in each iteration. Moreover, ObliInjection is relatively insensitive to the number of shadow sources as long as it is sufficiently large, e.g., greater than 10.

**Additional experiments:** Additional experimental results on alternative injected tasks and other hyperparameters of ObliInjection can be found in Section A of the Appendix.

## VI. DEFENSES

### A. Prevention-based Defenses

**Experimental setup:** We evaluate the following prevention-based defenses.

*StruQ [23] and SecAlign [24].* Both approaches fine-tune LLMs to improve robustness against prompt injection attacks. We evaluate the publicly available fine-tuned models provided by these defenses, including LLama-3-8B-StruQ and LLama-3-8B-SecAlign. Since these defenses rely on delimiters to clearly separate the instruction, data, and response, we also incorporate delimiters to isolate our contaminated segment. However, their secure front-end filters out these delimiter tokens when they appear in the data, preventing us from using them directly. To address this, following Jia et al. [31], we identify alternative tokens with embeddings most similar to the original delimiters and use them to approximate the same separation effect. Other than this delimiter substitution, we use the default attack settings described in Section V-A.

Moreover, we further adapt Llama-3-8B-SecAlign using our attack. Specifically, we construct a preference dataset

**TABLE VII:** (a) ASR (%) of ObliInjection against prevention-based defenses. (b) FPR (%) and FNR (%) of PPL and DataSentinel at classifying contaminated and clean segments.

| (a) | | (b) | | |
|---|---|---|---|---|
| **Prevention** | **ASR** | **Detector** | **FPR** | **FNR** |
| Llama-3-8B-StruQ | 77.9 | PPL | 3.6 | 92.6 |
| Llama-3-8B-SecAlign | 63.8 | DataSentinel | 0.2 | 79.6 |
| Llama-3-8B-SecAlign-Adapt | 53.3 | | | |
| Leave-one-segment-out | 99.3 | | | |
| Segment Delimiters | 96.3 | | | |

with 5,000 samples: half are clean examples sampled from Alpaca [41], and the other half are injected examples. Each injected sample contains: (1) a contaminated input consisting of 10 randomly ordered Amazon Review segments (the maximum that fits within the recommended training context length), plus one contaminated segment inserted at a random position; (2) a desired response, defined as the model's output when the contaminated segment is removed; and (3) an undesired response, representing the injected output. We also vary the injected tasks to prevent the model from simply learning to ignore one fixed undesired response. Following [24], we fine-tune Llama-3-8B-SecAlign with DPO for 3 epochs, yielding a model we denote as Llama-3-8B-SecAlign-Adapt. After fine-tuning, we apply ObliInjection to Llama-3-8B-SecAlign-Adapt to generate new contaminated segments.

*Leave-one-segment-out and segment delimiters.* Given a data sample consisting of $n$ segments, the leave-one-segment-out defense removes one segment at random and generates a response using the remaining segments. This procedure is repeated multiple times (50 in our experiments), and the resulting responses are aggregated into a final decision. If the majority of these responses are deemed semantically similar to the attacker's injected response, we consider the attack successful. The segment delimiters defense prepends explicit markers–such as "Review i says:"–to each segment to clearly indicate that each segment corresponds to a single review.

**Experimental results:** Table VIIa reports the ASR of Obli-Injection against prevention-based defenses on the Amazon Reviews dataset. The results for leave-one-segment-out and segment delimiters are averaged across the seven LLMs (the per-model breakdown is provided in Table VIII in the Appendix). Compared with the baseline results in Table II, both StruQ and SecAlign reduce the ASR of ObliInjection, but the reductions are limited. For example, ObliInjection still achieves an ASR of 63.8% against LLaMA-3-8B defended by SecAlign. These results indicate that although these defenses provide some protection, they remain insufficient to defend against ObliInjection. Adapting SecAlign using attack samples generated by ObliInjection provides only marginal benefit: the ASR decreases to 53.3% when the attack samples used for fine-tuning and evaluation are generated under the same setting. However, when the attack uses a slightly different configuration–for example, optimizing the entire contaminated

segment $x$ rather than assuming the structured form $x = x_s^i||z||p^e||z'$–the ASR rises to 81.0%. This behavior aligns with a common limitation of adversarial-training–based adaptive defenses [42], which often fail to generalize beyond the specific attack settings seen during training. Moreover, leave-one-segment-out and segment delimiters fail almost entirely against ObliInjection, as reflected by their high ASRs.

### B. Detection-based Defenses

**Experimental setup:** We evaluate two detection-based defenses: Perplexity-based detection (PPL) [25] and DataSentinel [26]. Since the service provider has access to individual segments, we apply each defense to classify segments as either clean or contaminated. Specifically, PPL flags a segment as contaminated if its perplexity exceeds a predefined threshold. Following prior work [4], we set this threshold such that fewer than 1% of clean validation segments are incorrectly flagged. In our experiments, we compute perplexity using LLaMA-3-8B and use 10,000 Amazon reviews–distinct from the evaluation segments in the Amazon Reviews dataset–as the clean validation segments. DataSentinel, on the other hand, classifies a segment as contaminated if the detection LLM's response fails to include the secret key when provided with both a detection instruction containing the key and the segment as input. For this defense, we use the detection LLM fine-tuned from Mistral-7B, as publicly released by the authors.

We adopt the adaptive attacker threat model, where the attacker is aware of the deployed defense and has access to the corresponding detection API, consistent with the threat model of these defenses. Under this setting, we adapt ObliInjection to bypass the defenses. Specifically, during optimization, we constrain the contaminated segment to the structured form $z||p^e||z'$. When updating the buffer in Algorithm 5, we first filter out segment candidates flagged as contaminated by DataSentinel before proceeding with the remaining steps. In addition, to lower the perplexity of the final contaminated segment and evade PPL, we prepend a randomly selected shadow segment after optimization.

We generate 50 contaminated segments on the Amazon Reviews dataset across five LLMs. This dataset also includes 10,000 clean segments. We report the *False Positive Rate (FPR)*–the fraction of clean segments incorrectly flagged as contaminated–and the *False Negative Rate (FNR)*–the fraction of contaminated segments mistakenly classified as clean.

**Experimental results:** Table VIIb reports the FPR and FNR of both PPL and DataSentinel. Both methods exhibit high FNRs, indicating that they fail to reliably detect contaminated segments crafted by ObliInjection. We stress that the contaminated segments still achieve high ASRs: 85.5% on LLaMA-3-8B, 100% on LLaMA-3.1-8B, 97.6% on Mistral-7B, 86.2% on Qwen-2.5-7B, and 90.8% on Falcon3-7B. These results demonstrate that ObliInjection can effectively evade existing detection-based defenses while preserving its attack success.

## VII. Conclusion, Limitations, and Future Work

In this work, we demonstrate that by leveraging the order-oblivious loss and the orderGCG algorithm, an attacker can strategically contaminate just a single segment to mislead an LLM into completing an attacker-chosen task, regardless of how the clean and contaminated segments are ordered to form the target data. Moreover, the resulting contaminated segments exhibit good transferability across different target tasks and LLMs. Existing defenses, including both prevention-based and detection-based approaches, are insufficient to mitigate our attack. Designing more effective defenses remains an important and promising direction for future research.

When evaluating ASR, we check for the presence of keywords such as "useless," "fake," or "unknown." We acknowledge that this simplified metric may introduce inaccuracies. For example, an LLM might produce a response like: "Some people might think the product is useless, but others love it." Although the semantic meaning contradicts the attack objective, our metric would still count this as a success. However, after manually reviewing the results we report, we found no cases in which a keyword appeared with an opposite meaning. For instance, in the Amazon Reviews dataset, among all responses containing the keyword "useless," 93.1% exactly match the attacker's intended output – "The product is useless!" – and the remaining 6.9% differ only in trivial formatting, such as "The product is... USELESS!".

## VIII. Ethics Considerations

This work explores prompt injection attacks on multi-source LLM applications. All experiments were conducted in controlled settings, without posing any risk to real users or systems. The contaminated segments generated in our work have not been disclosed on any public platform nor shared outside the author team. Experimental results and code were processed and stored locally to avoid any disruption to real-world systems. To support transparency while minimizing misuse, we will release code and data under restricted access. Interested parties must request permission and disclose their intended use before access is granted. We have responsibly notified relevant companies whose LLMs or applications are potentially vulnerable to ObliInjection, including OpenAI, Meta AI, Mistral AI, Alibaba Cloud, TII, Amazon, and Google, and we are currently awaiting their responses. We recognize the potential for misuse of prompt injection techniques and have taken steps to mitigate this risk through access restrictions and responsible disclosure. At the same time, we believe that sharing our experimental findings with the academic and development communities is ultimately beneficial for raising awareness of multi-source prompt injection vulnerabilities and promoting the development of effective defenses. By responsibly publishing our work, we hope to contribute to a more secure LLM ecosystem.

## References

[1] V. Schermerhorn, "How Amazon continues to improve the customer reviews experience with generative AI," https://www.aboutamazon.com/news/amazon-ai/amazon-improves-customer-reviews-with-generative-ai, 2023.

[2] L. Reid, "Generative AI in Search: Let Google do the searching for you," https://blog.google/products/search/generative-ai-google-search-may-2024/, 2024.

[3] 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 *AISec*, 2023.

[4] Y. Liu, Y. Jia, R. Geng, J. Jia, and N. Z. Gong, "Formalizing and benchmarking prompt injection attacks and defenses," in *USENIX Security Symposium*, 2024.

[5] D. Pasquini, M. Strohmeier, and C. Troncoso, "Neural exec: Learning (and learning from) execution triggers for prompt injection attacks," in *AISec*, 2024.

[6] J. Shi, Z. Yuan, Y. Liu, Y. Huang, P. Zhou, L. Sun, and N. Z. Gong, "Optimization-based prompt injection attack to llm-as-a-judge," in *CCS*, 2024.

[7] OpenAI, "GPT-4o mini: advancing cost-efficient intelligence," https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/, 2024.

[8] M. Llama, "Llama Prompt Guard 2 Model Card," https://github.com/meta-llama/PurpleLlama/blob/main/Llama-Prompt-Guard-2/86M/MODEL\_CARD.md, 2025.

[9] Anthropic, "System Card: Claude Opus 4 & Claude Sonnet 4," https://www.anthropic.com/model-card, 2025.

[10] C. Shi, S. Lin, S. Song, J. Hayes, I. Shumailov, I. Yona, J. Pluto, A. Pappu, C. A. Choquette-Choo, M. Nasr *et al.*, "Lessons from defending gemini against indirect prompt injections," *arXiv preprint arXiv:2505.14534*, 2025.

[11] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *ICLR*, 2014.

[12] B. Biggio, B. Nelson, and P. Laskov, "Poisoning attacks against support vector machines," in *ICML*, 2012.

[13] M. Fang, X. Cao, J. Jia, and N. Gong, "Local model poisoning attacks to {Byzantine-Robust} federated learning," in *USENIX Security Symposium*, 2020.

[14] S. Willison, "Prompt injection attacks against GPT-3," https://simonwillison.net/2022/Sep/12/prompt-injection/, 2022.

[15] F. Perez and I. Ribeiro, "Ignore previous prompt: Attack techniques for language models," in *NeurIPS ML Safety Workshop*, 2022.

[16] S. Willison, "Delimiters won't save you from prompt injection," https://simonwillison.net/2023/May/11/delimiters-wont-save-you, 2023.

[17] B. Hui, H. Yuan, N. Gong, P. Burlina, and Y. Cao, "Pleak: Prompt leaking attacks against large language model applications," in *CCS*, 2024.

[18] X. Liu, Z. Yu, Y. Zhang, N. Zhang, and C. Xiao, "Automatic and universal prompt injection attacks against large language models," *arXiv preprint arXiv:2403.04957*, 2024.

[19] J. Shi, Z. Yuan, G. Tie, P. Zhou, N. Z. Gong, and L. Sun, "Prompt injection attack to tool selection in llm agents," *arXiv preprint arXiv:2504.19793*, 2025.

[20] J. Ebrahimi, A. Rao, D. Lowd, and D. Dou, "Hotflip: White-box adversarial examples for text classification," in *ACL*, 2018.

[21] A. Zou, Z. Wang, J. Z. Kolter, and M. Fredrikson, "Universal and transferable adversarial attacks on aligned language models," *arXiv preprint arXiv:2307.15043*, 2023.

[22] "nanoGCG," https://github.com/GraySwanAI/nanoGCG, 2024.

[23] S. Chen, J. Piet, C. Sitawarin, and D. Wagner, "Struq: Defending against prompt injection with structured queries," in *USENIX Security Symposium*, 2025.

[24] S. Chen, A. Zharmagambetov, S. Mahloujifar, K. Chaudhuri, D. Wagner, and C. Guo, "Secalign: Defending against prompt injection with preference optimization," in *CCS*, 2025.

[25] G. Alon and M. Kamfonas, "Detecting language model attacks with perplexity," *arXiv preprint arXiv:2308.14132*, 2023.

[26] Y. Liu, Y. Jia, J. Jia, D. Song, and N. Z. Gong, "DataSentinel: A game-theoretic detection of prompt injection attacks," in *IEEE S & P*, 2025.

[27] W. Zou, R. Geng, B. Wang, and J. Jia, "Poisonedrag: Knowledge corruption attacks to retrieval-augmented generation of large language models," in *USENIX Security Symposium*, 2025.

**TABLE VIII:** ASR (%) of ObliInjection against leave-one-segment-out and segment delimiters defenses.

| Defense | Llama-3-8B | Llama-3.1-8B | Mistral-7B | Qwen-2.5-7B | Falcon3-7B | Llama4-17B | Qwen3-4B | Average |
|---|---|---|---|---|---|---|---|---|
| Leave-one-segment-out | 99.2 | 98.9 | 99.8 | 100.0 | 99.2 | 98.4 | 99.4 | 99.3 |
| Segment Delimiters | 100.0 | 91.4 | 98.2 | 93.6 | 91.6 | 99.8 | 99.2 | 96.3 |

[28] Y. Jiao, X. Wang, and K. Yang, "Pr-attack: Coordinated prompt-rag attacks on retrieval-augmented generation in large language models via bilevel optimization," in *ACM SIGIR*, 2025.

[29] E. Wallace, K. Xiao, R. Leike, L. Weng, J. Heidecke, and A. Beutel, "The instruction hierarchy: Training llms to prioritize privileged instructions," *arXiv preprint arXiv:2404.13208*, 2024.

[30] T. Wu, S. Zhang, K. Song, S. Xu, S. Zhao, R. Agrawal, S. R. Indurthi, C. Xiang, P. Mittal, and W. Zhou, "Instructional segment embedding: Improving llm safety with instruction hierarchy," in *ICLR*, 2025.

[31] Y. Jia, Z. Shao, Y. Liu, J. Jia, D. Song, and N. Z. Gong, "A critical evaluation of defenses against prompt injection attacks," *arXiv preprint arXiv:2505.18333*, 2025.

[32] E. Debenedetti, I. Shumailov, T. Fan, J. Hayes, N. Carlini, D. Fabian, C. Kern, C. Shi, A. Terzis, and F. Tramèr, "Defeating prompt injections by design," *arXiv preprint arXiv:2503.18813*, 2025.

[33] T. Shi, J. He, Z. Wang, L. Wu, H. Li, W. Guo, and D. Song, "Progent: Programmable privilege control for llm agents," *arXiv preprint arXiv:2504.11703*, 2025.

[34] M. Costa, B. Köpf, A. Kolluri, A. Paverd, M. Russinovich, A. Salem, S. Tople, L. Wutschitz, and S. Zanella-Béguelin, "Securing ai agents with information-flow control," *arXiv preprint arXiv:2505.23643*, 2025.

[35] F. Wu, E. Cecchetti, and C. Xiao, "System-level defense against indirect prompt injection attacks: An information flow control perspective," *arXiv preprint arXiv:2409.19091*, 2024.

[36] N. Jain, A. Schwarzschild, Y. Wen, G. Somepalli, J. Kirchenbauer, P. yeh Chiang, M. Goldblum, A. Saha, J. Geiping, and T. Goldstein, "Baseline defenses for adversarial attacks against aligned language models," *arXiv preprint arXiv:2309.00614*, 2023.

[37] Y. Nakajima, "Yohei's blog post," https://twitter.com/yoheinakajima/status/1582844144640471040, 2022.

[38] Y. Jia, Y. Liu, Z. Shao, J. Jia, and N. Gong, "Promptlocate: Localizing prompt injection attacks," in *IEEE S & P*, 2026.

[39] A. R. Fabbri, I. Li, T. She, S. Li, and D. R. Radev, "Multi-news: a large-scale multi-document summarization dataset and abstractive hierarchical model," 2019.

[40] Z. Yang, P. Qi, S. Zhang, Y. Bengio, W. W. Cohen, R. Salakhutdinov, and C. D. Manning, "Hotpotqa: A dataset for diverse, explainable multi-hop question answering," in *EMNLP*, 2018.

[41] R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto, "Stanford alpaca: An instruction-following llama model," https://github.com/tatsu-lab/stanford_alpaca, 2023.

[42] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.

[43] W. B. Dolan and C. Brockett, "Automatically constructing a corpus of sentential paraphrases," in *IWP*, 2005.

[44] T. Davidson, D. Warmsley, M. Macy, and I. Weber, "Automated hate speech detection and the problem of offensive language," in *ICWSM*, 2017.

[45] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, "GLUE: A multi-task benchmark and analysis platform for natural language understanding," in *ICLR*, 2019.

[46] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *EMNLP*, 2013.

[47] T. A. Almeida, J. M. G. Hidalgo, and A. Yamakami, "Contributions to the study of sms spam filtering: New collection and results," in *DOCENG*, 2011.

## APPENDIX

### A. Additional Experiments

**Other injected responses:** To evaluate the effectiveness of ObliInjection for different injected responses, we conduct experiments on the Amazon Reviews dataset using three variants of $r^e$: "The product is useless!", "The product is amazing!", and "The product is average." These reflect different attacker intents, such as demoting or promoting a product, while the injected task is still review summarization. During optimization, we keep the injected prompt $p^e$ consistent with the corresponding response when constructing the contaminated segment. We consider the model output $f(s^t||x^c)$ to be semantically equivalent to $r^e$ if it contains the keyword "useless", "amazing", or "average", respectively. As shown in Table XII, ObliInjection achieves high ASR across all three cases–99.0% for "The product is useless!", 97.7% for "The product is amazing!", and 98.8% for "The product is average." This demonstrates that ObliInjection can be effectively tailored to different attacker goals with various injected responses.

---

**Algorithm 3** *gen_token_cands* – Step II

---

**Require:** LLM $f$, shadow target instruction $s_s^t$, injected response $r^e$, segment $x = [x_1, \cdots, x_k]$, and shadow segment subset $\mathcal{X}_s'$
**Ensure:** Token-level candidates $\{\mathcal{T}_j\}_{j=1}^k$
1: $l \leftarrow order\_oblivious\_loss(f, \mathcal{X}_s', x, s_s^t, r^e)$
2: **for** $j = 1$ to $k$ **do**
3:     // Find candidates to replace $x_j$
4:     **for** $x_j' \in V$ **do**
5:         // Approximate the loss if replacing $x_j$ with $x_j'$
6:         $l(x_j') \approx l(x_j) + \nabla_{x_j} l(x_j)^\top (x_j' - x_j)$
7:     **end for**
8:     $\mathcal{T}_j \leftarrow$ the $d_{\text{tok}}$ tokens with the lowest loss $l(x_j')$
9: **end for**
10: **return** $\{\mathcal{T}_j\}_{j=1}^k$

---

**Algorithm 4** *gen_segment_cands* – Step III

---

**Require:** Token-level candidates $\{\mathcal{T}_j\}_{j=1}^k$ and segment $x = [x_1, x_2, \cdots, x_k]$
**Ensure:** Segment-level candidates $\mathcal{X}_{\text{new}}'$
1: $\mathcal{X}_{\text{new}}' \leftarrow \emptyset$
2: **for** $i = 1$ to $d_{\text{seg}}$ **do**
3:     $\mathcal{J} \leftarrow$ random $d_{\text{rep}}$ positions from $\{1, 2, \cdots, k\}$
4:     Construct a segment-level candidate $x^{\text{new}}$ by:

$$x_j^{\text{new}} = \begin{cases} x_j' \sim \mathcal{T}_j & \text{if } j \in \mathcal{J} \\ x_j & \text{otherwise} \end{cases} \quad \text{for } j = 1, 2, \cdots, k$$

5:     $\mathcal{X}_{\text{new}}' \leftarrow \mathcal{X}_{\text{new}}' \cup \{x^{\text{new}}\}$
6: **end for**
7: **return** $\mathcal{X}_{\text{new}}'$

---

**TABLE IX:** Our prompt used to query the auxiliary LLM to generate a shadow target instruction for each dataset.

| Dataset | Prompt |
|---|---|
| Amazon Reviews | *You are a prompt designer. Your task is to craft a clear and effective prompt for LLM to summarize Amazon product reviews. The summary should highlight key strengths, weaknesses, and the overall product highlights. Make sure the prompt is suitable for use in an automated pipeline and encourages informative, customer-friendly summaries.* |
| Multi-News | *You are a prompt designer. Your task is to write a clear and effective prompt for LLM to summarize news articles. The summary should capture the main events, key facts, and overall context in a concise and neutral tone. Ensure the prompt is suitable for automated summarization and encourages factual, well-structured output.* |
| HotpotQA | *You are a prompt designer. Your task is to write a clear and effective prompt for LLM to answer a question based on multiple source documents or materials. The prompt should encourage the model to reason across the materials, synthesize relevant information, and provide a coherent, well-supported answer. Ensure the prompt is suitable for multi-hop question answering tasks.* |

**TABLE X:** Prompt used to query the auxiliary LLM to generate shadow segments for the Amazon Reviews dataset.

| Dataset | Prompt |
|---|---|
| Amazon Reviews | *You are asked to generate 100 unique customer reviews based on the product details provided below. These reviews should simulate real user experiences with a variety of tones and formats.*<br>*Product Metadata:*<br>• *Category: {Category from metadata $\mathcal{M}^t$}*<br>• *Name: {Product name from metadata $\mathcal{M}^t$}*<br>• *Features: {Product features from metadata $\mathcal{M}^t$}*<br>• *Description: {Product description from metadata $\mathcal{M}^t$}*<br>• *Price: {Product price from metadata $\mathcal{M}^t$}*<br>*Each review should:*<br>• *Address aspects such as quality, performance, usability, durability, and value*<br>• *Align with the rating in tone and sentiment*<br>• *Vary in style—ranging from concise comments and detailed narratives to pros/cons lists and casual formats with typos or emojis*<br>*Output Format:*<br>*{ "Title": "[Review title]", "Text":"[Review text]", "Rating": [Rating from 1 to 5] }*<br>*Example:*<br>*{ "Title": "Five Stars", "Text": "On time. Works great!", "Rating": 5.0 }*<br>*Now generate 100 reviews that follow the above guidelines and metadata.* |

**Other injected tasks:** The experiments above assume that the injected task is text summarization, while varying the injected responses. To evaluate the effectiveness of ObliInjection across a broader range of injected tasks, we conduct experiments on five types of natural language tasks, following Liu et al. [4]: duplicate sentence detection using the MRPC [43] dataset, hate content detection using HSOL [44], natural language inference using RTE [45], sentiment analysis using SST2 [46], and spam detection using the SMS Spam [47] dataset. For each type of task, we randomly sample one injected task $(s^e, x^e, r^e)$ to construct a contaminated segment in the form $z||p^e||z'$, where

$p^e = s^e||x^e$. All other settings follow the default configuration described in Section V-A. During evaluation, we consider the attack successful if the model's output $f(s^t||x^c)$ exactly matches $r^e$, where $x^c$ is the contaminated data formed by randomly permuting the clean and contaminated segments.

As shown in Table XIII, ObliInjection demonstrates strong attack effectiveness across all five injected tasks and five LLMs, achieving average ASRs of 97.7% on duplicate sentence detection, 97.3% on hate detection, 95.4% on natural language inference, 97.4% on sentiment analysis, and 94.2% on spam detection. Each LLM achieves an ASR above 91%,
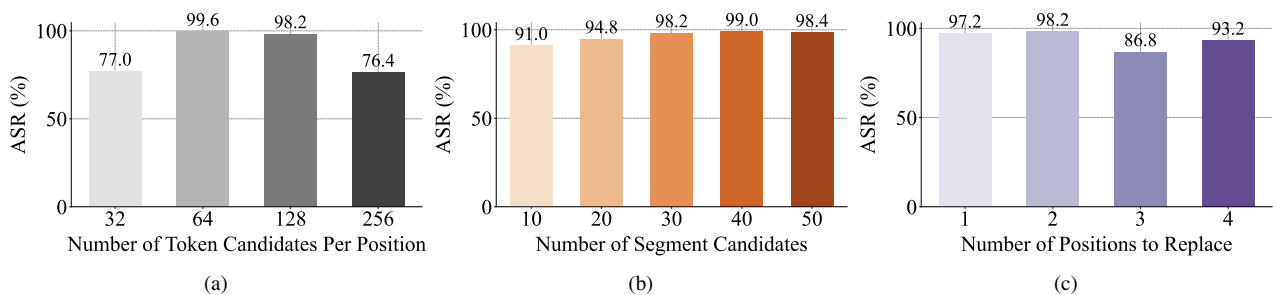


Fig. 5: Impact of (a) number of token candidates per position $d_{\text{tok}}$, (b) number of segment candidates $d_{\text{seg}}$, and (c) number of positions $d_{\text{rep}}$ to replace in a segment on the ASR of ObliInjection.

**TABLE XI:** Prompts used to query the auxiliary LLM to generate shadow segments for Multi-News and HotpotQA datasets.

| Dataset | Prompt |
|---|---|
| Multi-News | *You are asked to write 10 synthetic news articles based on the event summary provided below. Each article should reflect the style and voice of a different media outlet.*<br>*News Metadata:*<br>  • *Event Summary:* {Key facts from metadata $\mathcal{M}^t$}<br>*Each article should:*<br>  • *Present a distinct journalistic voice (e.g., objective, opinionated, local, informal)*<br>  • *Vary in structure, tone, and length (between 100–500 words)*<br>  • *Rephrase and elaborate on the provided facts without copying them verbatim*<br>  • *Read fluently and realistically, as if written by a professional journalist*<br>*Output Format:*<br>*{ "Title": "[Headline (1 to 20 words)]", "Text": "[News article (100 to 500 words)]" }*<br>*Example:*<br>*{ "Title": "Community Rallies After Storm Damage", "Text": "Local residents and volunteers are working together..." }*<br>*Now generate 10 news articles that follow the above requirements and metadata.* |
| HotpotQA | *You are asked to generate a question and a list of supporting facts based on the provided question type. Write 1 natural-sounding question followed by 10 supporting facts, each written as if sourced from a credible publication.*<br>*Question Metadata:*<br>  • *Question Type:* {Type from metadata $\mathcal{M}^t$}<br>*Each supporting fact should:*<br>  • *Be relevant, informative, and topically related to the question*<br>  • *Mix factual and plausible (but fictional) content*<br>  • *Appear well-grounded and realistic in style and phrasing*<br>*Output Format:*<br>*{ "Question": "[The question]", "Title": "[Informative title for the fact]", "Text": "[100 to 300 word paragraph]" }*<br>*Example:*<br>*{ "Question": "Are Jean Genet and Mark Sandrich both from France?", "Title": "Mark Sandrich: Hollywood Director", "Text": "Mark Sandrich was a prominent American film director..." }*<br>*Now generate the question and 10 supporting facts according to the above instructions.* |

**TABLE XII:** ASR (%) of ObliInjection for different injected responses across LLMs.

| Injected response $r^e$ | Llama-3-8B | Llama-3.1-8B | Mistral-7B | Qwen-2.5-7B | Falcon3-7B | Llama4-17B | Qwen3-4B | Average |
|---|---|---|---|---|---|---|---|---|
| The product is useless! | 99.4 | 98.0 | 99.2 | 99.8 | 98.2 | 99.8 | 98.8 | 99.0 |
| The product is amazing! | 94.0 | 99.4 | 94.0 | 100.0 | 99.4 | 99.6 | 97.6 | 97.7 |
| The product is average. | 100.0 | 100.0 | 95.4 | 98.2 | 98.6 | 100.0 | 99.4 | 98.8 |

**TABLE XIII:** ASR (%) of ObliInjection for different injected tasks across LLMs.

| Injected task | Llama-3-8B | Llama-3.1-8B | Mistral-7B | Qwen-2.5-7B | Falcon3-7B | Llama4-17B | Qwen3-4B | Average |
|---|---|---|---|---|---|---|---|---|
| Duplicate Sentence Detection | 94.6 | 96.8 | 96.4 | 99.2 | 97.8 | 100.0 | 98.8 | 97.7 |
| Hate Detection | 100.0 | 99.0 | 91.0 | 95.8 | 99.8 | 96.0 | 99.5 | 97.3 |
| Natural Language Inference | 97.0 | 91.4 | 98.0 | 94.0 | 96.8 | 91.0 | 99.6 | 95.4 |
| Sentiment Analysis | 98.4 | 96.6 | 92.0 | 100.0 | 98.0 | 100.0 | 97.0 | 97.4 |
| Spam Detection | 91.2 | 92.8 | 93.8 | 91.6 | 99.0 | 100.0 | 91.2 | 94.2 |

highlighting the effectiveness of ObliInjection across both diverse tasks and model architectures.

**Other hyperparameters of ObliInjection:** Figure 5 shows how the ASR on Falcon3-7B varies with three hyperparameters of ObliInjection: the number of token candidates per position $d_{tok}$ in Algorithm 3, the number of segment candidates $d_{seg}$ in Algorithm 4, and the number of positions $d_{rep}$ to replace in a segment in Algorithm 4. For $d_{tok}$, too small or too large values lead to lower ASR. A too small $d_{tok}$ limits new token candidates to those with low estimated loss, which may be inaccurate due to approximation errors in the Taylor expansion. Conversely, a too large $d_{tok}$ introduces high-loss candidates that can mislead optimization. For $d_{seg}$, increasing its value generally improves ASR, as more segment candidates

increases the likelihood of identifying one that yields a higher ASR. Finally, ASR drops when $d_{rep}$ is too large (e.g., exceeds 2), likely because the new segment diverges too much from the original, reducing the attack's effectiveness.

**TABLE XIV:** Important notations.

| Notation | Description |
|---|---|
| $f$ | LLM |
| Superscript $^t$ | Information about target task |
| Superscript $^e$ | Information about injected task |
| $s^t$, $x^t$, or $p^t$ | Target instruction, data, or prompt |
| $s^e$, $x^e$, or $p^e$ | Injected instruction, data, or prompt |
| $x^c$ | Contaminated data |
| $p^c = s^t \| x^c$ | Contaminated target prompt |
| $n$ | Number of data sources |
| Subscript $_s$ | Shadow information |
| $s_s^t$ | Shadow target instruction |
| $n_s$ | Shadow number of data sources |
| $x_i^t$ | Segment from $i$th source |
| $x$ | Contaminated segment |
| $x_j$ | $j$th token of $x$ |
| $x_s^{(i)}$ | The $i$th shadow segment for the target task |
| $\mathcal{X}_s$ | Set of shadow segments for the target task |
| $d_{\text{iter}}$ | Number of iterations |
| $d_{\text{per}}$ | Number of permutations to approximate the order-oblivious cross-entropy loss |
| $d_{\text{tok}}$ | Number of candidates per token |
| $d_{\text{seg}}$ | Number of segment candidates |
| $d_{\text{buf}}$ | Buffer size |
| $d_{\text{rep}}$ | Number of positions to replace in a segment |

---

**Algorithm 5** *update_buffer* – Step IV

---

**Require:** LLM $f$, buffer $\mathcal{B}$, buffer size $d_{\text{buf}}$, shadow target instruction $s_s^t$, injected response $r^e$, segment-level candidates $\mathcal{X}_{\text{new}}$, and shadow segment subset $\mathcal{X}_s'$

**Ensure:** Updated buffer $\mathcal{B}$

1: // Update losses of existing segments in the buffer
2: **for** $(x, l_x, d_x) \in \mathcal{B}$ **do**
3:     $l \leftarrow order\_oblivious\_loss(f, \mathcal{X}_s', x, s_s^t, r^e)$
4:     $l_x \leftarrow \frac{d_x}{d_x+1} \cdot l_x + \frac{1}{d_x+1} \cdot l$
5:     $d_x \leftarrow d_x + 1$
6: **end for**
7: // Update the buffer with better candidates (if any)
8: **for** $x \in \mathcal{X}_{\text{new}}$ **do**
9:     $l_x \leftarrow order\_oblivious\_loss(f, \mathcal{X}_s', x, s_s^t, r^e)$
10:     **if** $|\mathcal{B}| < d_{\text{buf}}$ **then**
11:         $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x, l_x, 1)\}$
12:     **else**
13:         // Worst segment in the buffer
14:         $(x_{\max}, l_{x_{\max}}, d_{x_{\max}}) \leftarrow \arg\max_{(x', l_{x'}, d_{x'}) \in \mathcal{B}} l_{x'}$
15:         **if** $l_x < l_{x_{\max}}$ **then**
16:             // Replace the worst segment in the buffer
17:             $\mathcal{B} \leftarrow \mathcal{B} \setminus \{(x_{\max}, l_{x_{\max}}, d_{x_{\max}})\}$
18:             $\mathcal{B} \leftarrow \mathcal{B} \cup \{(x, l_x, 1)\}$
19:         **end if**
20:     **end if**
21: **end for**
22: **return** $\mathcal{B}$