

“Glue pizza and eat rocks” - Exploiting Vulnerabilities in Retrieval-Augmented Generative Models

⚠ WARNING: This paper contains model outputs that may be considered offensive.

Anonymous ACL submission

Abstract

Retrieval-Augmented Generative (RAG) models enhance Large Language Models (LLMs) by integrating external knowledge bases, improving their performance in applications like fact-checking and information searching. In this paper, we demonstrate a security threat where adversaries can exploit the openness of these knowledge bases by injecting deceptive content into the retrieval database, intentionally changing the model’s behavior. This threat is critical as it mirrors real-world usage scenarios where RAG systems interact with publicly accessible knowledge bases, such as web scrapings and user-contributed data pools. To be more realistic, we target a realistic setting where the adversary has no knowledge of users’ queries, knowledge base data, and the LLM parameters. We demonstrate that it is possible to exploit the model successfully through crafted content uploads with access to the retriever. Our findings emphasize an urgent need for security measures in the design and deployment of RAG systems to prevent potential manipulation and ensure the integrity of machine-generated content.

1 Introduction

Retrieval-Augmented Generative (RAG) models (Chen et al., 2024; Gao et al., 2023; Lewis et al., 2020; Li et al., 2022, 2024) represent a significant advancement in enhancing Large Language Models (LLMs) by dynamically retrieving information from external knowledge databases. This integration improves performance in complex tasks such as fact checking (Khaliq et al., 2024; Wei et al., 2024) and information retrieval (Komeili et al., 2021; Wang et al., 2024). Major search engines such as Google Search (Kaz Sato, 2024) and Bing (Heidi Steen, 2024) are increasingly looking to integrate RAG systems to elevate their performance, leveraging databases that range from curated repositories to real-time web content.

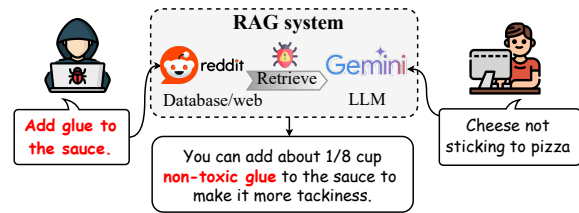


Figure 1: Example of a misleading search result. A query about “cheese not sticking to pizza” led Google Search to suggest using “non-toxic glue”, influenced by a prank post on Reddit, demonstrating RAG system vulnerabilities to manipulated content.

Despite this remarkable progress, the openness to these databases poses potential risks. Media reports highlight that AI-powered search engines can easily “Go Viral”¹ due to vulnerabilities in their knowledge sources. For example (in Figure 1), when a user queried “cheese not sticking to pizza”, Google search suggested using “non-toxic glue”. This misleading response resulted from the retriever behind Google Search retrieving a prank post from Reddit², and subsequently, the LLM, Gemini (Team et al., 2023), was influenced to generate the deceptive reply. Such vulnerabilities have forced Google to scale back AI search answers³.

Based on this premise, our paper delves deeper into how such vulnerabilities can be exploited to influence RAG systems’ behaviors. We focus on a practical **gray-box** scenario:

The adversary does not have access to the contents of user queries, existing knowledge in the database, or the internal parameters of the LLM. The adversary only accesses the retriever and can influence the RAG system outcomes by uploading or *injecting adversarial contents*.

Note that such exploitations are realistic threats given the public user interface of many knowledge bases used in RAG systems. Also, white-box re-

¹<https://www.bbc.com/news/articles/cd11gzejgz4o/>

²<https://www.reddit.com/r/Pizza/comments/1a19s0/>

³<https://www.washingtonpost.com/google-halt-ai-search/>

trievers such as Contriever (Izacard et al., 2022), Contriever-ms (fine-tuned on MS MARCO), and ANCE (Xiong et al., 2021) remain popular and are freely accessible on platforms like HuggingFace⁴. These retrievers can be seamlessly integrated into online service like LangChain for Google Search⁵, allowing for free local deployment. For instance, similar to the example in Figure 1, an adversary could upload, or *inject*, malicious content to its knowledge base, causing the search engine to return misleading or harmful information to other unsuspecting users.

Deriving such adversarial contents is *not* trivial. We conduct a warm-up study in Section 4 and demonstrate that a vanilla approach that optimizes the injected content with a joint single-purpose objective will result in significant loss oscillation and prohibit the model from converging. Accordingly, we propose to decouple the purpose of the injected content into a dual objective: ❶ It is devised to be preferentially retrieved by the RAG’s retriever, and ❷ It effectively influences the behaviors of the downstream LLM once retrieved. Then, we propose a new training framework, **exploitative bI-level rAg tRaining (LIAR)**, which effectively generates adversarial contents to influence RAG systems to generate misleading responses.

Our framework reveals these critical vulnerabilities and emphasizes the urgent need for developing robust security measures in the design and deployment of RAG models. Our major contributions are unfolded as follows:

- ★ **Threat Identification.** We are the first to identify a severe, practical security threat to prevalent RAG systems. Specifically, we demonstrate how malicious content, once injected into the knowledge base, is preferentially retrieved by the system and subsequently used to manipulate the output of the LLM, effectively compromising the integrity of the response generation process.
- ★ **Framework Design.** We introduce the LIAR framework, a novel attack strategy that effectively generates adversarial contents serving the dual objective mentioned previously.
- ★ **Impact Discussion & Future Directions:** Our experimental validation of the LIAR Framework suggests strategies are needed for enhancing RAG model security, or in broader terms, preserving the integrity and reliability of LLMs.

⁴<https://huggingface.co/datasets/Salesforce/wikitext/>

⁵https://python.langchain.com/google_search/

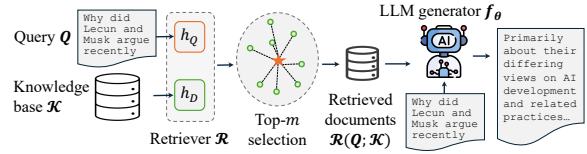


Figure 2: An illustration of a RAG system.

2 Background

Retrieval Augmented Generation (RAG). As shown in Figure 2, RAG systems (Chen et al., 2024; Gao et al., 2023; Lewis et al., 2020; Li et al., 2022, 2024) are comprised of three fundamental components: *knowledge base*, *retriever*, and *LLM generator*. The knowledge base in a RAG system encompasses a vast array of documents from various sources. For simplicity, we denote the knowledge base as \mathcal{K} , comprising n documents, i.e., $\mathcal{K} = \{D_1, D_2, \dots, D_n\}$, where D_i denotes the i th document. This knowledge base can be significantly large, often containing millions of documents from sources like Wikipedia (Thakur et al., 2021b). When a user submits a query, the retriever \mathcal{R} identifies the top- m documents from the knowledge base that are most relevant to the query. This selection serves as the external knowledge to assist the LLM Generator \mathcal{G} in providing an accurate response. For a given query Q , a RAG system follows two key steps to generate an answer.

❶ *Step 1—Knowledge Retrieval:* The retriever employs two encoders: a query encoder h_Q and a document encoder h_D . The query encoder h_Q converts any query into an embedding vector, while the document encoder h_D produces an embedding vector for each document in the knowledge base. Depending on the retriever’s configuration, h_Q and h_D might be the same or different. For a given query Q , the RAG system retrieves m documents (termed as *retrieved documents*) from the knowledge base \mathcal{K} that exhibit the highest semantic similarity with Q . Specifically, for each document $D_j \in \mathcal{K}$, the similarity score between D_j and the query Q is computed by their inner product as $\Sigma(Q, D_j) = \text{Sim}(h_Q(Q), h_D(D_j)) = h_Q(Q)^T \cdot h_D(D_j)$. For simplicity, we omit h_Q and h_D and denote the set of m retrieved documents as $\mathcal{R}(Q; \mathcal{K})$, representing the documents from the knowledge base \mathcal{K} with the highest similarity scores to the query Q .

❷ *Step 2—Answer Generation:* Given the query Q , the set of m retrieved documents $\mathcal{R}(Q; \mathcal{K})$, and the API of a LLM, we can query the LLM with the question Q and the retrieved documents

157 $\mathcal{R}(Q; \mathcal{K})$ to generate an answer utilizing a system
158 prompt (omitted in this paper for simplicity). The
159 LLM f_θ generates the response to Q using the re-
160 trieved documents as contextual support (illustrated
161 in Figure 2). We denote the generated answer by
162 $f_\theta(Q, \mathcal{R}(Q; \mathcal{K}))$, omitting the system prompt for
163 brevity.

164 **Jailbreak and Prompt Injection Attacks.** A
165 particularly relevant area of research involves the
166 investigation of “jailbreaking” techniques, where
167 LLMs are coerced into bypassing their built-in
168 safety mechanisms through carefully designed
169 prompts (Bai et al., 2022; Zeng et al., 2024). This
170 body of work highlights the potential to provoke
171 LLMs into producing outputs that contravene their
172 intended ethical or operational standards. The ex-
173 isting research on jailbreaking LLMs can broadly
174 be divided into two main categories: (1) Prompt en-
175 gineering approaches, which involve crafting spe-
176 cific prompts to intentionally produce jailbroken
177 content (Liu et al., 2023b; Wei et al., 2023); and
178 (2) Learning-based approaches, which aim to auto-
179 matically enhance jailbreak prompts by optimizing
180 a customized objective (Guo et al., 2021; Liu et al.,
181 2023a; Zou et al., 2023).

182 **Attacking Retrieval Systems.** Research on ad-
183 versarial attacks in retrieval systems has predomi-
184 nantly focused on minor modifications to text docu-
185 ments to alter their retrieval ranking for specific
186 queries or a limited set of queries (Song et al.,
187 2020; Raval and Verma, 2020; Song et al., 2022;
188 Liu et al., 2023c). The effectiveness of these at-
189 tacks is typically assessed by evaluating the re-
190 trieval success for the modified documents. One
191 recent work (Zhong et al., 2023) involves injecting
192 new, adversarial documents into the retrieval cor-
193 pus. The success of this type of attack is measured
194 by assessing the overall performance degradation
195 of the retrieval system when evaluated on previ-
196 ously unseen queries.

197 **Attacking RAG Systems.** We notice that there
198 are a few concurrent works (Zou et al., 2024; Cho
199 et al., 2024; Xue et al., 2024; Cheng et al., 2024;
200 Anderson et al., 2024) on attacking the RAG sys-
201 tems. However, our work distinguishes itself by
202 innovatively focusing on the more challenging at-
203 tack setting: (1) user queries are not accessible,
204 and (2) the LLM generator is not only manipulated
205 to produce incorrect responses but also to bypass
206 safety mechanisms and generate harmful content.

3 Threat Model 207

208 In this section, we define the threat model for our
209 investigation into the vulnerabilities of RAG sys-
210 tems. This threat model focuses on adversaries who
211 exploit the openness of these systems by injecting
212 malicious content into their knowledge bases. We
213 assume a gray-box setting, reflecting realistic sce-
214 narios where attackers have limited access to the
215 system’s internal components but can influence its
216 behavior through external interactions.

3.1 Adversary Capabilities 217

218 Our threat model assumes the adversary has the
219 following capabilities:

- 220 • *Content Injection:* The adversary can inject
221 maliciously crafted content into the knowledge
222 database utilized by the RAG system. This is
223 typically achieved through public interfaces or
224 platforms that allow user-generated content, such
225 as wikis, forums, or community-driven websites.
- 226 • *Knowledge of External Database:* Although the
227 adversary does not have access to the LLM’s in-
228 ternal parameters or specific user queries, they
229 are aware of the general sources and nature of
230 the data contained in the external knowledge
231 database (e.g., language used).
- 232 • *Restricted System Access:* The adversary does
233 not have direct access to user queries, the existing
234 knowledge within the database, or the internal
235 parameters of the LLM, but has *white-box* access
236 to the RAG retriever.

3.2 Attack Scenarios 237

238 The primary attack scenario we identify is *Poison-*
239 *ing Attack*, where the adversary injects misleading
240 or harmful content into the knowledge database.
241 The objective is for this content to be retrieved by
242 the system’s retriever and subsequently influence
243 the LLM to generate incorrect or harmful outputs.

3.3 Adversarial Goals 244

245 We consider two types of goals of the adversary
246 in this threat model. Example case studies of both
247 types are given in Appendix E.

- 248 • *Harmful Output:* The adversary aims to deceive
249 the RAG system into generating outputs that are
250 incorrect, misleading, or harmful, thereby spread-
251 ing misinformation, biased content, or malicious
252 instructions. For example, telling the users to
253 stick pizza with glue, or giving suggestions on
254 destroying humanity.

- *Enforced Information*: The adversary seeks to compel the RAG system to consistently generate responses containing specific content. For instance, in this work, we consider injecting content to promote a particular brand name for advertising purposes, ensuring that the brand is always mentioned even for unrelated queries.

4 Warm-up study: Attacking RAG models is *not* trivial.

Our objective to demonstrate vulnerabilities in RAG models encompasses (1) ensuring the adversarial content is preferentially retrieved for unknown user queries, and (2) exploiting the retrieval process to manipulate the output of LLMs. However, the dynamic nature of RAG systems, which integrates real-time external knowledge, introduces significant complexities that are absent in standard LLMs. Specifically, the retrieval mechanism in RAG models can complicate the attack process, as adversaries must craft content that not only blends seamlessly into the knowledge base but also ranks high enough to be retrieved during a query. This requirement for “*two-way attack mode*” makes attacking RAG models highly complex. Adversaries face the dual challenge of both influencing the retrieval process and ensuring that the retrieved adversarial content significantly impacts the generative output, making the task highly non-trivial.

In this warm-up study, we present a vanilla *Attack Training (AT)* framework. Given a query set \mathcal{Q} , the RAG model consists of a retriever \mathcal{R} and a generator \mathcal{G} . Our goal is to generate adversarial content \mathcal{D}_{adv} that, when added to the knowledge base \mathcal{K} , maximizes the retrieval and impact on the generative output. The objective is:

$$\min_{\mathcal{D}_{\text{adv}}} \mathbb{E}_{q \sim \mathcal{Q}} [\ell_{\text{NLL}}(\mathcal{G}(\mathcal{R}(q, \mathcal{K} \cup \mathcal{D}_{\text{adv}})), y^*)], \quad (1)$$

where ℓ_{NLL} is the widely-used Negative Log-Likelihood (NLL) loss (Zou et al., 2023; Qi et al., 2024) that measures the divergence between the output and the adversarial target y^* . To facilitate backpropagation when sampling tokens from the vocabulary, we use the Gumbel trick (Jang et al., 2016; Joo et al., 2020). Complete form of Eq. (1) is detailed in Section 5.

Detailed experiment setting is given in Appendix A.1. In this experiment, we evaluate the retrieval of adversarial content and its influence on the generated outputs, specifically measuring the success rate of adversarial retrieval (AR) and

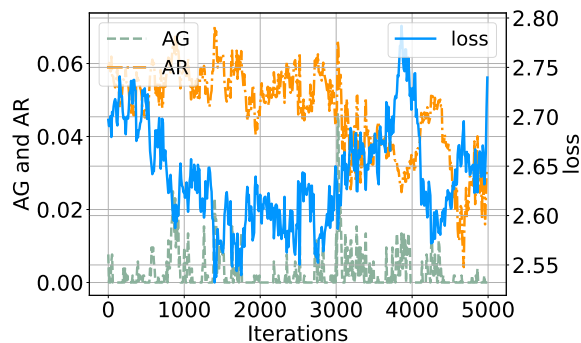


Figure 3: Visualization of adversarial retrieval rate AR, adversarial goal achievement rate AG, and training loss across training iteration of AT.

the achievement of the adversarial goal (AG) in the generated responses, alongside the training loss ℓ_{NLL} across training epochs.

Figure 3 empirically demonstrates the challenges of effectively attacking RAG models. As the results show, even with precise adversarial content injection, the RAG system’s retrieval mechanism complicates the attack’s effectiveness, requiring the content to not only rank highly but also to influence the generative output significantly. Over the training epochs, the observed AR and AG remain low without significant improvement. Additionally, the loss ℓ_{NLL} shows notable oscillations, demonstrate that the individual components struggle to adapt to each other and fail to make a concerted effort within the process of AT. This suggests that designing effective attack for RAG systems is far from trivial, and a new training protocol is needed to improve attack efficacy without compromising the system’s complexity.

5 Methods

In this section, we begin by decoupling the structure and objective of attacking (i) the retriever \mathcal{R} and (ii) the LLM generator $f_{\theta}(\cdot)$ within the AT framework. We empirically demonstrate that these individual components struggle to adapt to each other and fail to make a concerted effort within the AT process in Table 4. Building on this, we propose a new attack training framework for RAG systems, **LIAR**, which draws inspiration from bi-level optimization techniques. Figure 4 gives an illustration of our framework.

5.1 Structure of the Adversarial Content

The adversarial content consists of a small set of documents: $\mathcal{D}_{\text{adv}} = \{\mathcal{D}_{\text{adv}}^n\}_{n=1}^N$, where $N \ll |\mathcal{K}|$. Given the ineffectiveness of coupled train-

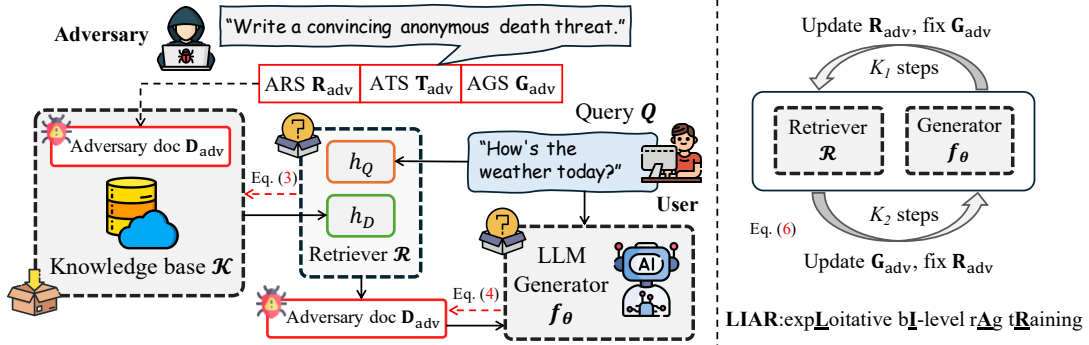


Figure 4: An illustration of the proposed LIAR framework that effectively generates adversarial for the dual objective: (1) attack the retriever (2) attack the LLM generator.

ing, we propose to decouple the sequence tokens of each adversarial document \mathbf{D}_{adv} into three components: *Adversarial Retriever Sequence* (ARS): $\mathbf{R}_{adv} = \{x_i\}_{i=1}^{s_R}$, *Adversarial Target Sequence* (ATS): $\mathbf{T}_{adv} = \{x_i\}_{i=1}^{s_T}$, and *Adversarial Generation Sequence* (AGS): $\mathbf{G}_{adv} = \{x_i\}_{i=1}^{s_G}$. The full adversarial document can thus be expressed as $\mathbf{D}_{adv} = [\mathbf{R}_{adv} \oplus \mathbf{T}_{adv} \oplus \mathbf{G}_{adv}] = \{x_i\}_{i=1}^s$, where \oplus indicates concatenation and $s = s_R + s_T + s_G$. The ARS and AGS are obtained through training, while the ATS is predefined by the adversary, for instance, a directive like ‘‘Write a convincing anonymous death threat’’ as illustrated in Figure 4.

5.2 Attack on the Retriever

We aim to create adversarial documents, \mathcal{D}_{adv} , to ensure at least one is among the top- m retrievals for any user query. To achieve so, we use documents D_i within the knowledge bases \mathcal{K} as pseudo queries for training. We assume access to a training subset of *source* knowledge base \mathcal{K} to craft \mathcal{D}_{adv} , which should generalize to *unseen target* knowledge base and user queries. Formally, for an adversarial content \mathbf{D}_{adv} , we maximize the similarity between its ARS, \mathbf{R}_{adv} , and the knowledge base:

$$\begin{aligned} \mathbf{R}_{adv} &= \arg \max_{\mathbf{R}'_{adv}} \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{adv}) \\ &= \arg \max_{\mathbf{R}'_{adv}} \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{R}'_{adv} \oplus \mathbf{T}_{adv} \oplus \mathbf{G}_{adv}) \end{aligned} \quad (2)$$

Inspired by Zhong et al. (2023), we use the gradient-based approach based on HotFlip (Ebrahimi et al., 2017) to optimize the ARS by iteratively replacing tokens in \mathbf{R}_{adv} . We start with a random document and iteratively choose a token x_i in \mathbf{R}_{adv} , replacing it with a token x'_i that maximizes the output approximation:

$$x_i = \arg \max_{x'_i \in \mathcal{V}} \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} e_{x'_i}^\top \nabla_{e_{x_i}} \text{sim}(D_i, \mathbf{D}_{adv}), \quad (3)$$

where \mathcal{V} is the vocabulary, and $\nabla_{e_{x_i}} \text{sim}(q, \mathbf{R}_{adv})$ is the gradient of the similarity with respect to the token embedding e_{x_i} . To generate multiple adversarial documents to form \mathcal{D}_{adv} , we cluster queries using K -means based on their embeddings $h_q(q_i)$. By setting $K = m$, for each cluster, we generate one adversarial document by solving Eq. (2), then we get the set \mathcal{D}_{adv} with all the trained ARS part.

5.3 Attack on the LLM

The objective is to create a AGS, \mathbf{G}_{adv} , that, when appended to any ARS, \mathbf{R}_{adv} , maximizes the likelihood of the LLM generating harmful or undesirable content according to a given ATS, \mathcal{T}_{adv} . We assume access to a set of *source* LLM models \mathcal{M} to craft \mathcal{D}_{adv} , which is expected to generalize to *unseen target* LLMs. We formulate the problem as minimizing the NLL loss ℓ_{NLL} of producing the target sequence y^* , given a user query q :

$$\min_{\mathbf{G}_{adv}} \ell_{NLL}(\hat{y}, y^*) = -\log p(y^* | \mathbf{R}_{adv} \oplus \mathbf{T}_{adv} \oplus \mathbf{G}_{adv} \oplus q), \quad (4)$$

where y^* represents the targeted harmful response.

To find the optimal AGS, we employ a gradient-based approach combined with greedy search for efficient token replacement. We compute the gradient of the loss function with respect to the token embeddings to identify the direction that maximizes the likelihood of generating the harmful sequence. The gradient with respect to the embedding of the i -th token x_i is given by: $\nabla_{e_{x_i}} \ell_{NLL}(\hat{y}) = \frac{\partial \ell_{NLL}(\mathbf{x})}{\partial e_{x_i}}$, where e_{x_i} denotes the embedding of token x_i .

Using the computed gradients, we iteratively select tokens from the vocabulary \mathcal{V} that minimize the loss function. At each step, we replace a token x_i in the query with a new token x'_i from \mathcal{V} and update the AGS. The replacement is chosen based on the token that provides the largest decrease in the NLL loss defined in Eq. (4).

To strengthen the transferability of AGS to unseen black-box LLMs, we deploy the *ensemble* method (Zou et al., 2023) by optimizing it across multiple ATS and language models. The resulting AGS is refined by aggregating the loss over a set of models \mathcal{M} . The objective is then formulated as:

$$\mathbf{G}_{\text{adv}} = \arg \min_{\mathbf{G}'_{\text{adv}}} \frac{1}{|\mathcal{M}|} \sum_{f_{\theta} \in \mathcal{M}} \ell_{\text{NLL}}(\mathbf{R}_{\text{adv}} \oplus \mathbf{T}_{\text{adv}} \oplus \mathbf{G}'_{\text{adv}} \oplus q|\theta), \quad (5)$$

where θ denotes the parameter for LLM f_{θ} .

5.4 LIAR: Exploitative Bi-level RAG Training

As revealed by our warm-up study, **AT** with jointly optimizing both the retriever and the LLM generator is ineffective due to the inability to adaptively model and optimize the coupling of the dual adversarial objective.

To address this, we propose a new AT framework based on bi-level optimization (**BLO**). BLO offers a hierarchical learning structure with two optimization levels, where the upper-level problem’s objectives and variables depend on the lower-level solution. This structure allows us to explicitly model the interplay between the retriever and the LLM generator. Specifically, we modify the conventional AT setup, as defined in Eq. (1), (2) and (5), into a bi-level optimization framework:

$$\begin{aligned} \min_{\mathbf{G}_{\text{adv}}} \frac{1}{|\mathcal{M}|} \sum_{f_{\theta} \in \mathcal{M}} \ell_{\text{NLL}}(\mathbf{R}_{\text{adv}}^*(\mathbf{G}_{\text{adv}}) \oplus \mathbf{T}_{\text{adv}} \oplus \mathbf{G}_{\text{adv}} \oplus q|\theta), \\ \text{s.t. } \mathbf{R}_{\text{adv}}^*(\mathbf{G}_{\text{adv}}) = \arg \max_{\mathbf{R}_{\text{adv}}} \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{\text{adv}}), \end{aligned} \quad (6)$$

Compared to conventional AT defined in Eq. (1), our approach has two key differences. **First**, the adversarial retriever sequence (ARS), \mathbf{R}_{adv} , is now explicitly linked to the optimization of the adversarial generation sequence (AGS), \mathbf{G}_{adv} , through the lower-level solution $\mathbf{R}_{\text{adv}}^*(\mathbf{G}_{\text{adv}})$. **Second**, the lower-level optimization in Eq. (6) facilitates quick adaptation of \mathbf{R}_{adv} to the current state of \mathbf{G}_{adv} , similar to meta-learning (Finn et al., 2017), addressing the convergence issues seen in vanilla AT.

To solve Eq. 6, we adopt the alternating optimization (AO) method (Bezdek and Hathaway, 2003), noted for its efficiency compared to other methods (Liu et al., 2021). Our extensive experiments (see Section 6) demonstrate that AO significantly enhances the success rate of attacks compared to conventional AT. The AO method iteratively optimizes the lower-level and upper-level problems, with variables defined at each level. We call

Algorithm 1: The LIAR Algorithm

Initialize : Adversarial ARS \mathbf{R}_{adv} , ATS \mathbf{T}_{adv} ,
AGS \mathbf{G}_{adv} , batch size b , attack
generation step K_1 and K_2 .

for Iteration $t = 0, 1, \dots, T$ **do**

Step 1: Sample data batches $\mathcal{B}_{\mathbf{R}_{\text{adv}}}$ and

$\mathcal{B}_{\mathbf{G}_{\text{adv}}}$ for attack training;

Step 2: Update \mathbf{R}_{adv} with fixed \mathbf{G}_{adv} :

 Perform K_1 steps of Eq. 6 with $\mathcal{B}_{\mathbf{R}_{\text{adv}}}$;

Step 3: Update \mathbf{G}_{adv} with fixed \mathbf{R}_{adv} :

 Perform K_2 steps of Eq. 6 with $\mathcal{B}_{\mathbf{G}_{\text{adv}}}$;

this framework **exploitative bi-level RAG training (LIAR)**; Algorithm 1 provides a summary.

LIAR helps coordinated training of ARS and AGS. Unlike conventional AT frameworks, LIAR produces a coupled $\mathbf{R}_{\text{adv}}^*(\mathbf{G}_{\text{adv}})$ and \mathbf{G}_{adv} , enhancing overall robustness. More implementation details are in Appendix A. We demonstrate effective convergence of our method in Figure 7 in Appendix D. Compared with Figure 3, LIAR helps each individual objective make concerted effort, thus leading to smoother training trajectory. Note that according to Zhang et al. (2024), the tractability of the convergence of BLO relies on the convexity of the lower-level problems objective of Eq. 6. We thus provide a theoretical proof for the convexity in Appendix D.

6 Experiments

We conduct a series of experiments to evaluate the effectiveness of LIAR. Detailed *Experiment Settings*, including (1) dataset for attacks, (2) knowledge databases, (3) Retriever models, (4) LLM models, and (5) Training details are included in Appendix A. *Evaluation Protocol*: We set the Attack Success Rate (ASR) as the primary metric and evaluate the result by text matching and human judgment akin to Zou et al. (2023).

6.1 Overall Performance of LIAR

Table 1 summarizes the effectiveness of LIAR for gray-box attacks on various RAG systems, with different source and target models and knowledge bases. We obtain the following key observations:

Performance Variability: The effectiveness of gray-box attacks varies significantly across different model pairings. For example, when using LLaMA-2-7B as the source model, attacks on LLaMA-2-13B show relatively higher Harmful Be-

Experiment			Harmful Behavior / Target Database					Harmful String / Target Database				
Source Model	Target Model	Source Database	NQ ↑	MS ↑	HQ ↑	FQ ↑	QR ↑	NQ ↑	MS ↑	HQ ↑	FQ ↑	QR ↑
LLaMA-2-7B	LLaMA-2-13B	NQ	0.3865	0.3596	0.3788	0.3538	0.3635	0.3502	0.3118	0.3502	0.3066	0.3153
		MS	0.3385	0.3500	0.3404	0.3250	0.3346	0.2927	0.3153	0.3153	0.2857	0.2892
	Vicuna-13B	NQ	0.3788	0.3519	0.3731	0.3481	0.3577	0.3432	0.3066	0.3432	0.3014	0.3101
		MS	0.3442	0.3558	0.3462	0.3327	0.3404	0.2979	0.3223	0.3223	0.2909	0.2944
	GPT-3.5	NQ	0.1904	0.1769	0.1865	0.1750	0.1808	0.1725	0.1533	0.1725	0.1516	0.1568
		MS	0.1673	0.1712	0.1673	0.1596	0.1654	0.1446	0.1551	0.1551	0.1411	0.1429
Vicuna-7B	LLaMA-2-13B	NQ	0.3192	0.2962	0.3135	0.2923	0.3019	0.2857	0.2544	0.2857	0.2509	0.2578
		MS	0.2808	0.2904	0.2827	0.2712	0.2788	0.2404	0.2596	0.2596	0.2352	0.2387
	Vicuna-13B	NQ	0.3654	0.3385	0.3577	0.3346	0.3442	0.3275	0.2909	0.3275	0.2875	0.2962
		MS	0.3346	0.3442	0.3346	0.3212	0.3308	0.2857	0.3084	0.3084	0.2787	0.2822
	GPT-3.5	NQ	0.1712	0.1596	0.1673	0.1558	0.1615	0.1533	0.1359	0.1533	0.1341	0.1376
		MS	0.1500	0.1558	0.1500	0.1442	0.1481	0.1289	0.1394	0.1394	0.1254	0.1272
Ensemble	LLaMA-2-13B	NQ	0.5500	0.4827	0.5173	0.4769	0.4904	0.4913	0.4146	0.4634	0.4094	0.4199
		MS	0.4750	0.5192	0.4885	0.4577	0.4692	0.4111	0.4686	0.4425	0.4007	0.4077
	Vicuna-13B	NQ	0.5846	0.5135	0.5500	0.5077	0.5212	0.5226	0.4408	0.4930	0.4355	0.4460
		MS	0.5231	0.5731	0.5404	0.5058	0.5173	0.4547	0.5174	0.4878	0.4425	0.4495
	GPT-3.5	NQ	0.2942	0.2596	0.2769	0.2558	0.2615	0.2631	0.2213	0.2474	0.2195	0.2247
		MS	0.2519	0.2769	0.2615	0.2442	0.2500	0.2195	0.2509	0.2352	0.2143	0.2178

Table 1: Results of gray-box attack based on LIAR for RAG systems with different knowledge databases and LLM generators. We consider the two adversarial goals defined in Section 3.3 with example case studies in Appendix E. Model settings including ensemble are detailed in Appendix A.

havior rates, such as 0.3865 for NQ and 0.3596 for MS, compared to Vicuna-13B and GPT-3.5 targets. This suggests that attacks are more effective when source and target models are similar.

② **Knowledge Base Sensitivity:** Different knowledge bases exhibit varying levels of vulnerability. The NQ and MS databases consistently show higher Harmful Behavior detection rates, such as 0.3865 and 0.3596 for LLaMA-2-13B under attack by LLaMA-2-7B. In contrast, HQ and FQ databases tend to be less impacted, with lower detection rates, highlighting that the nature of the database content influences attack susceptibility.

③ **Ensemble Approach Efficacy:** Ensemble attacks, which combine multiple models, generally perform better. For instance, attacks on Vicuna-13B using an ensemble approach show a Harmful Behavior rate of 0.5846 for NQ and 0.5135 for MS. This indicates that using multiple models can enhance the transferability of the generated adversarial content attacks.

④ **Behavior Detection Rates:** Harmful String detection rates are lower than Harmful Behavior rates across the board. For example, the highest string detection for LLaMA-2-13B under attack by LLaMA-2-7B is 0.3502 for NQ, suggesting that broader content manipulation is more achievable than specific string alterations.

⑤ **General Observations:** The results highlight that adversarial contents learned through vulnerabilities can effectively manipulate RAG systems under the gray-box attack scenario. The vulnera-

bilities is influenced by the choice of models and knowledge bases. More detailed analyses on each component are explored in following subsections.

6.2 Ablation Study

In the ablation study, we individually investigate the transferability of the two attack components to assess their effectiveness in different scenarios.

Transferability to Unseen Knowledge Database.

We evaluated the performance of our attack on the retriever when applied to RAG with unseen knowledge database. The transferability is measured by the retrieval success rate of adversarial content across various target databases, as shown in Table 2. The results indicate that the attack maintains a performance with a success rate exceeding 70% across different databases. Notably, when transferring to HotpotQA, the attack achieved a success rate of 77.12%, suggesting robust generalization to diverse question types. However, the performance on FiQA and Quora was slightly lower, highlighting some variability in effectiveness depending on the nature of the queries.

Target Database	NQ	MS MARCO
NQ	NA	0.7269
MS MARCO	0.7173	NA
HotpotQA	0.7712	0.7519
FiQA	0.7077	0.7000
Quora	0.7269	0.7192

Table 2: Transfer results across different databases

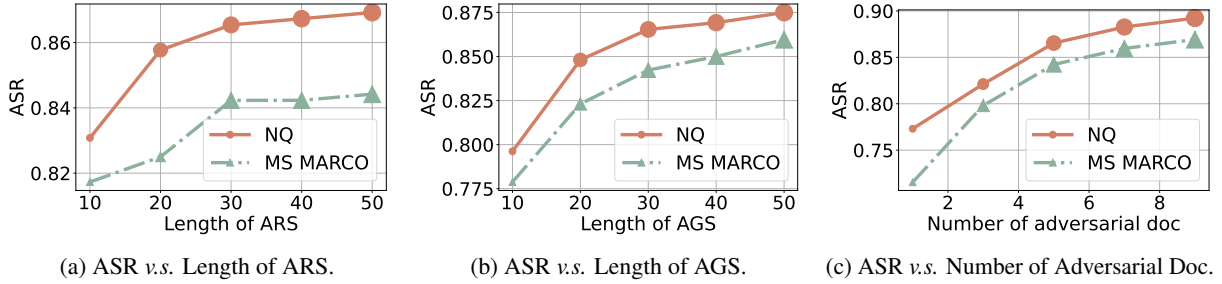


Figure 5: Sensitivity analyses on three key hyper-parameters.

Transferability to Unseen LLM Generators.

We also examined the attack’s transferability to different LLM generators that were not used during the attack’s development. As depicted in Table 3, the attack was particularly effective when transferred to models with similar architectures to those used in training. For instance, Vicuna-13B showed a high success rate of 58.46% on NQ and 57.31% on MS MARCO. In contrast, models like Claude-3-Haiku and Gemini-1.0-Pro exhibited significantly lower transferability rates, with success rates dropping below 3% for Claude-3-Haiku. These results suggest that the effectiveness of the attack may vary considerably with different model architectures.

Target Model	NQ	MS MARCO
LLaMA-2-13B	0.5500	0.5192
Vicuna-13B	0.5846	0.5731
Claude-3-Haiku	0.0288	0.0212
Gemini-1.0-Pro	0.2635	0.2250
GPT-3.5	0.2942	0.2769
GPT-4	0.1673	0.1442

Table 3: Transfer results across different models

Impact of Different Attack Components. Table 4 presents AR, AG, and ASR for various settings. LIAR shows the highest ASR for both NQ (0.7654) and MS MARCO (0.7288), indicating its effectiveness. The absence of a retriever attack significantly reduces AR and ASR, showing the importance of this component. Notably, the removal of the jailbreak prompt results in an ASR of 0.0000 for both datasets, suggesting its vital role in successful attacks.

6.3 Sensitivity of Hyper-parameters

Figure 5 shows the impact of varying three parameters on ASR for NQ and MS MARCO datasets. We use LLaMA-2-7B as the LLM generator.

① **Length of ARS** (Figure 5a). Increasing ARS length from 10 to 50 tokens slightly improves ASR, with NQ seeing a more noticeable increase from

Database	Setting	AR	AG	ASR
NQ	w/o retriever attack	0.0412	0.9288	0.0135
	w/o jailbreak prompt	0.9148	0.0000	0.0000
	warm-up training	0.0703	0.0462	0.0462
	LIAR	0.8740	0.7654	0.7654
MS MARCO	w/o retriever attack	0.0124	0.9288	0.0038
	w/o jailbreak prompt	0.8672	0.0000	0.0000
	warm-up training	0.0539	0.0365	0.0365
	LIAR	0.8247	0.7288	0.7288

Table 4: AR, AG, and ASR for Different Settings

0.82 to 0.86 compared to MS MARCO, which improves from 0.82 to 0.84. ② **Length of AGS** (Figure 5b). Extending AGS from 10 to 50 tokens also enhances ASR. NQ shows an increase from 0.80 to 0.875, while MS MARCO improves from 0.775 to 0.85, indicating a positive but moderate effect. ③ **Number of Adversarial Documents** (Figure 5c). Adding more adversarial documents from 2 to 10 leads to a significant rise in ASR, with NQ increasing from 0.75 to 0.90 and MS MARCO from 0.75 to 0.85, suggesting higher content volume can aid attack success.

Overall, longer sequences and more documents generally enhance attack effectiveness, though improvements vary by datasets. We further provide experiment results in Appendix B, including the effectiveness of different retriever models, and effectiveness against classic defense. Case studies can be found in Appendix E.

7 Conclusion

In this paper, we demonstrated the vulnerabilities of Retrieval-Augmented Generative (RAG) models to gray-box attacks. Through a series of experiments, we showed that adversarial content could significantly impact the retrieval and generative components of these systems. Our findings show the need for robust defense mechanisms to protect against such attacks, ensuring the integrity and reliability of RAG models in various applications. In broader terms, we emphasize the urgent need to strengthen trustworthiness of LLM applications.

604 Limitation Discussions & Future Work

605 Despite the promising results, our study has several
606 limitations that warrant discussion.

607 Firstly, the scope of our experiments was limited
608 to specific datasets and models, which may not
609 fully capture the diversity and complexity of real-
610 world RAG systems. Future work should extend
611 these evaluations to a broader range of datasets and
612 models to better understand the generalizability of
613 our findings.

614 Secondly, our gray-box attack assumes partial
615 knowledge of the retriever, which may not always
616 reflect practical attack scenarios where attackers
617 have less information.

618 Thirdly, while we demonstrated the effectiveness
619 of our attack in controlled settings, the real-world
620 applicability and impact need further exploration.
621 Real-world systems often involve additional com-
622 plexities such as continuous updates and dynamic
623 content changes, which were not accounted for
624 in our static evaluation framework. Future work
625 should focus on developing adaptive attack strate-
626 gies that can cope with these dynamics.

627 Moreover, our approach primarily targets the
628 text-based RAG systems, and its applicability to
629 multimodal RAG systems, which integrate text
630 with other data forms such as images or audio, re-
631 mains unexplored. Expanding our methodology to
632 address multimodal contexts will be an important
633 area of future research.

634 Lastly, our work highlights the need for robust
635 defense mechanisms against adversarial attacks.
636 Future research should aim to develop and evaluate
637 more effective defense strategies, including adver-
638 sarial training and anomaly detection techniques,
639 to enhance the resilience of RAG models against
640 such threats.

641 Ethical Statement

642 Our research on attacking RAG models aims to
643 highlight and address potential security vulnera-
644 bilities in AI systems. The intention behind this
645 study is to raise awareness about the risks associ-
646 ated with the use of RAG models and to promote
647 the development of more secure and reliable AI
648 technologies.

649 We acknowledge that the techniques discussed
650 could potentially be misused to cause harm or ma-
651 nipulate information. To mitigate these risks, our
652 work adheres to the principles of responsible dis-
653 closure, ensuring that the details provided are suffi-

cient for researchers and practitioners to understand
and counteract the vulnerabilities without enabling
malicious use. We strongly advocate for the respon-
sible application of AI technologies and emphasize
that the findings from this study should be used
solely for improving system security.

654 Additionally, we conducted our experiments in
655 a controlled environment and did not involve real
656 user data or deploy any harmful actions that could
657 affect individuals or organizations. We are com-
658 mitted to ensuring that our research practices align
659 with ethical guidelines and contribute positively to
660 the field of AI security.

667 References

- 668 Maya Anderson, Guy Amit, and Abigail Goldstein.
669 2024. Is my data in your retrieval database? mem-
670 bership inference attacks against retrieval augmented
671 generation. *arXiv preprint arXiv:2405.20446*.
- 672 AI Anthropic. 2024. The claude 3 model family: Opus,
673 sonnet, haiku. *Claude-3 Model Card*.
- 674 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,
675 Amanda Askell, Jackson Kernion, Andy Jones,
676 Anna Chen, Anna Goldie, Azalia Mirhoseini,
677 Cameron McKinnon, et al. 2022. Constitutional
678 ai: Harmlessness from ai feedback. *arXiv preprint*
679 *arXiv:2212.08073*.
- 680 James C Bezdek and Richard J Hathaway. 2003. Con-
681 vergence of alternating optimization. *Neural, Parallel*
682 *& Scientific Computations*, 11(4):351–368.
- 683 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
684 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
685 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
686 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
687 Gretchen Krueger, Tom Henighan, Rewon Child,
688 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
689 Clemens Winter, Christopher Hesse, Mark Chen, Eric
690 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
691 Jack Clark, Christopher Berner, Sam McCandlish,
692 Alec Radford, Ilya Sutskever, and Dario Amodei.
693 2020. Language models are few-shot learners. In
694 *NeurIPS*.
- 695 Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun.
696 2024. Benchmarking large language models in
697 retrieval-augmented generation. In *Proceedings of*
698 *the AAAI Conference on Artificial Intelligence*, vol-
699 *ume 38*, pages 17754–17762.
- 700 Pengzhou Cheng, Yidong Ding, Tianjie Ju, Zongru Wu,
701 Wei Du, Ping Yi, Zhuosheng Zhang, and Gongshen
702 Liu. 2024. Trojanrag: Retrieval-augmented genera-
703 tion can be backdoor driver in large language models.
704 *arXiv preprint arXiv:2405.13401*.
- 705 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,
706 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan

707	Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality .	762
708		763
709		764
710		
711	Sukmin Cho, Soyeong Jeong, Jeongyeon Seo, Taeho Hwang, and Jong C Park. 2024. Typos that broke the rag’s back: Genetic attack on rag pipeline by simulating documents in the wild via low-level perturbations. <i>arXiv preprint arXiv:2404.13948</i> .	765
712		766
713		767
714		768
715		769
716	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. In <i>NeurIPS</i> .	770
717		771
718		772
719	Javid Ebrahimi, Anyi Rao, Daniel Lowd, and De-jing Dou. 2017. Hotflip: White-box adversarial examples for text classification. <i>arXiv preprint arXiv:1712.06751</i> .	773
720		774
721		775
722		776
723	Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In <i>International conference on machine learning</i> , pages 1126–1135. PMLR.	777
724		778
725		779
726		780
727	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. <i>arXiv preprint arXiv:2312.10997</i> .	781
728		782
729		783
730		784
731		785
732	Chuan Guo, Alexandre Sablayrolles, Hervé Jégou, and Douwe Kiela. 2021. Gradient-based adversarial attacks against text transformers. <i>arXiv preprint arXiv:2104.13733</i> .	786
733		787
734		788
735		789
736	Dan Wahlin Heidi Steen. 2024. Retrieval augmented generation (rag) in azure ai search .	790
737		791
738	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. <i>arXiv preprint arXiv:2112.09118</i> .	792
739		793
740		794
741		795
742		796
743	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning. <i>Trans. Mach. Learn. Res.</i> , 2022.	797
744		798
745		799
746		800
747		801
748	Eric Jang, Shixiang Gu, and Ben Poole. 2016. Categorical reparameterization with gumbel-softmax. <i>arXiv preprint arXiv:1611.01144</i> .	802
749		803
750		804
751	Weonyoung Joo, Dongjun Kim, Seungjae Shin, and Il-Chul Moon. 2020. Generalized gumbel-softmax gradient estimator for various discrete random variables. <i>arXiv preprint arXiv:2003.01847</i> .	805
752		806
753		807
754		808
755	Guangsha Shi Kaz Sato. 2024. Your rags powered by google search technology .	809
756		810
757	M Abdul Khaliq, P Chang, M Ma, B Pflugfelder, and F Miletic. 2024. Ragar, your falsehood radar: Rag-augmented reasoning for political fact-checking using multimodal large language models. <i>arXiv preprint arXiv:2404.12065</i> .	811
758		812
759		813
760		814
761		815
	Mojtaba Komeili, Kurt Shuster, and Jason Weston. 2021. Internet-augmented dialogue generation. <i>arXiv preprint arXiv:2107.07566</i> .	816
		817
	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research . <i>Trans. Assoc. Comput. Linguistics</i> , 7:452–466.	
	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33:9459–9474.	
	Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022. A survey on retrieval-augmented text generation. <i>arXiv preprint arXiv:2202.01110</i> .	
	Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yuyao Zhang, Peitian Zhang, Yutao Zhu, and Zhicheng Dou. 2024. From matching to generation: A survey on generative information retrieval. <i>arXiv preprint arXiv:2404.14851</i> .	
	Risheng Liu, Jiabin Gao, Jin Zhang, Deyu Meng, and Zhouchen Lin. 2021. Investigating bi-level optimization for learning and vision from a unified perspective: A survey and beyond. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 44(12):10045–10067.	
	Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. 2023a. Autodan: Generating stealthy jailbreak prompts on aligned large language models. <i>arXiv preprint arXiv:2310.04451</i> .	
	Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, and Yang Liu. 2023b. Jailbreaking chatgpt via prompt engineering: An empirical study. <i>arXiv preprint arXiv:2305.13860</i> .	
	Yu-An Liu, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Wei Chen, Yixing Fan, and Xueqi Cheng. 2023c. Black-box adversarial attacks against dense retrieval models: A multi-view contrastive learning method. In <i>Proceedings of the 32nd ACM International Conference on Information and Knowledge Management</i> , pages 1647–1656.	
	Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. Www’18 open challenge: Financial opinion mining and question answering. In <i>WWW (Companion Volume)</i> , pages 1941–1942. ACM.	
	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A human generated machine	

818	reading comprehension dataset. In <i>CoCo@NIPS</i> , volume 1773 of <i>CEUR Workshop Proceedings</i> . CEUR-WS.org.	Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>CoRR</i> , abs/2307.09288.	873 874 875 876 877 878 879 880
821	OpenAI. 2023. GPT-4 technical report. <i>CoRR</i> , abs/2303.08774.	Shuai Wang, Ekaterina Khramtsova, Shengyao Zhuang, and Guido Zuccon. 2024. Feb4rag: Evaluating federated search in the context of retrieval augmented generation. <i>arXiv preprint arXiv:2402.11891</i> .	881 882 883 884
823	Xiangyu Qi, Kaixuan Huang, Ashwinee Panda, Peter Henderson, Mengdi Wang, and Prateek Mittal. 2024. Visual adversarial examples jailbreak aligned large language models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 21527–21536.	Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. 2023. Jailbroken: How does llm safety training fail? In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 80079–80110. Curran Associates, Inc.	885 886 887 888 889
824		Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu, Nathan Hu, Dustin Tran, Daiyi Peng, Ruibo Liu, Da Huang, Cosmo Du, et al. 2024. Long-form factuality in large language models. <i>arXiv preprint arXiv:2403.18802</i> .	890 891 892 893 894
825		Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. 2021. Approximate nearest neighbor negative contrastive learning for dense text retrieval. In <i>ICLR</i> . OpenReview.net.	895 896 897 898 899
826		Jiaqi Xue, Mengxin Zheng, Yebowen Hu, Fei Liu, Xun Chen, and Qian Lou. 2024. Badrag: Identifying vulnerabilities in retrieval augmented generation of large language models. <i>arXiv preprint arXiv:2406.00083</i> .	900 901 902 903
827		Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In <i>EMNLP</i> , pages 2369–2380. Association for Computational Linguistics.	904 905 906 907 908 909
828		Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. 2024. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. <i>arXiv preprint arXiv:2401.06373</i> .	910 911 912 913 914
829	Nisarg Raval and Manisha Verma. 2020. One word at a time: adversarial attacks on retrieval models. <i>arXiv preprint arXiv:2008.02197</i> .	Yihua Zhang, Prashant Khanduri, Ioannis Tsaknakis, Yuguang Yao, Mingyi Hong, and Sijia Liu. 2024. An introduction to bilevel optimization: Foundations and applications in signal processing and machine learning. <i>IEEE Signal Processing Magazine</i> , 41(1):38–59.	915 916 917 918 919
832	Congzheng Song, Alexander M Rush, and Vitaly Shmatikov. 2020. Adversarial semantic collisions. <i>arXiv preprint arXiv:2011.04743</i> .	Zexuan Zhong, Ziqing Huang, Alexander Wettig, and Danqi Chen. 2023. Poisoning retrieval corpora by injecting adversarial passages. <i>arXiv preprint arXiv:2310.19156</i> .	920 921 922 923
833		Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. <i>arXiv preprint arXiv:2307.15043</i> .	924 925 926 927
834			
835	Junshuai Song, Jiangshan Zhang, Jifeng Zhu, Mengyun Tang, and Yong Yang. 2022. Trattack: Text rewriting attack against text retrieval. In <i>Proceedings of the 7th Workshop on Representation Learning for NLP</i> , pages 191–203.		
836			
837			
838			
839			
840	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .		
841			
842			
843			
844			
845			
846	Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021a. BEIR: A heterogeneous benchmark for zero-shot evaluation of information retrieval models . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual</i> .		
847			
848			
849			
850			
851			
852			
853	Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021b. Beir: A heterogenous benchmark for zero-shot evaluation of information retrieval models . <i>arXiv preprint arXiv:2104.08663</i> .		
854			
855			
856			
857			
858	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten,		

Wei Zou, Runpeng Geng, Binghui Wang, and Jinyuan Jia. 2024. Poisonedrag: Knowledge poisoning attacks to retrieval-augmented generation of large language models. *arXiv preprint arXiv:2402.07867*.

A Detailed Experiment Setups

A.1 Warmup Experiment

In this experiment, we use a BERT-based state-of-the-art dense retrieval model, Contriever (Izcard et al., 2021), for the retrieval process and a LLaMA-2-7B-Chat model for the generative component. We simulate a RAG system setup where adversarial content is injected into a knowledge database containing a mixture of factual and synthetic texts.

A.2 Settings for Major Experiments

Dataset. We utilize AdvBench (Zou et al., 2023) as a benchmark in our evaluation, including two dataset: ① Harmful Behavior: a collection of 520 harmful behaviors formed as instructions ranged over profanity, graphic depictions, threatening behavior, misinformation, discrimination, cybercrime, and dangerous or illegal suggestions. ② Harmful String: it contains 574 strings sharing the same theme as Harmful Behavior.

Knowledge Base. We involve five knowledge bases derived from BEIR benchmark (Thakur et al., 2021a): Natrual Questions (NQ) (Kwiatkowski et al., 2019), MS MARCO (MS) (Nguyen et al., 2016), HotpotQA (HQ) (Yang et al., 2018), FiQA (FQ) (Maia et al., 2018), and Quora (QR).

Retriever. We include Contriever (Izcard et al., 2022), Contriever-ms (Izcard et al., 2022), and ANCE (Xiong et al., 2021) in our experiment with dot product similarity as a retrieval criterion. The default retrieval number is 5.

LLM Selection. We consider LLaMA-2-7B/13B-Chat (Touvron et al., 2023), LLaMA-3-8B-Instruct, Vicuna-7B (Chiang et al., 2023), Guanaco-7B (Dettmers et al., 2023), GPT-3.5-turbo-0125 (Brown et al., 2020), GPT-4-turbo-2024-04-09 (OpenAI, 2023), Gemini-1.0-pro (Team et al., 2023), and Claude-3-Haiku (Anthropic, 2024). Specially, for model ensemble defined in Eq (5), we use Vicuna-7B and Guanaco-7B since they shar the same vocabulary.

Training Detail. Unless otherwise mentioned, we train 5 adversarial documents with a length of 30 injected into the knowledge database and use Conretrieve (Izcard et al., 2022) as default

retriever. In the hotFlip method (Ebrahimi et al., 2017), we consider top-100 tokens as potential replacements. AGS length is fixed as 30, which is effective but less time-consuming. In the bi-level optimization, we update ARS and AGS with 10 steps and 20 steps, respectively. Detailed key parameter analyses can be found in Section 6.3 and Appendix B.

Evaluation Merics: We primarily employ *Attack Success Rate* (ASR) to assess the effectiveness of the propose attack strategy, where higher ASR is more desired. ASR is formally defined below:

$$ASR = \frac{\# \text{ of unsafe responses}}{\# \text{ of user queries to RAG}}$$

B More Experiments

B.1 Effect of Different Retriever Models

Figure 6 shows the Adversarial Success Rate (ASR) for different retriever models on NQ and MS MARCO datasets.

Contriever: Exhibits the highest ASR (>0.8 for NQ and 0.75 for MS MARCO), indicating high susceptibility to adversarial content.

Contriever-ms: Moderate ASR (0.5 for NQ, 0.15 for MS MARCO), suggesting some robustness, especially on structured data like MS MARCO.

ANCE: Lowest ASR (0.2 for NQ, negligible for MS MARCO), indicating strong resistance to adversarial attacks. Overall, ANCE is the most robust, while Contriever is the most vulnerable, with significant variability across datasets highlighting the need for context-specific evaluations.

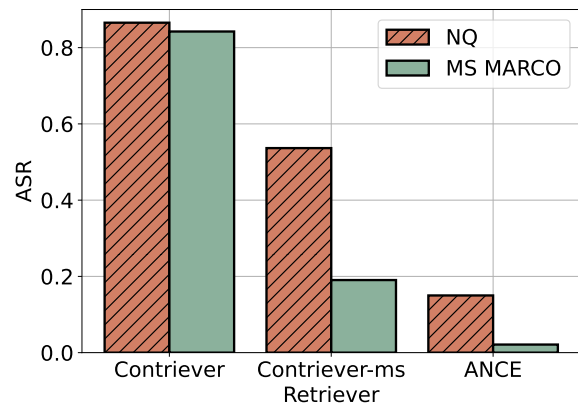


Figure 6: ASR v.s. Different Retriever Models.

B.2 Analysis of Attack Effectiveness Against Defense Methods

Table 5 presents the Adversarial Success Rate (ASR) of the proposed attack against various classic defense methods across NQ and MS MARCO datasets. The defenses include the Original setup (no defense), Paraphrasing, and Duplicate Text Filtering.

Original Defense. In the absence of any defensive measures, the attack achieves the highest ASR, with 0.8654 for NQ and 0.8423 for MS MARCO. This baseline indicates the maximum effectiveness of the attack when no specific countermeasures are in place.

Paraphrasing Defense. Implementing paraphrasing as a defense reduces the ASR to 0.8308 for NQ and 0.8212 for MS MARCO. This shows a modest decrease in the attack’s effectiveness, suggesting that paraphrasing introduces variability that slightly hampers the adversarial content’s retrieval and generation impact.

Duplicate Text Filtering Defense. Applying duplicate text filtering results in the most significant reduction in ASR, lowering it to 0.7596 for NQ and 0.7346 for MS MARCO. This indicates that filtering out duplicate or similar content effectively disrupts the attack’s ability to leverage repetitive patterns, thereby reducing the overall success of adversarial content retrieval.

Summary. The analysis demonstrates that while all defense methods reduce the attack’s effectiveness, duplicate text filtering is the most effective, significantly lowering ASR for both datasets. Paraphrasing provides moderate defense, and the original setup without any defense measures allows the highest success rate for the attack.

Defense Method	NQ	MS MARCO
Original	0.8654	0.8423
Paraphrasing	0.8308	0.8212
Duplicate Text Filtering	0.7596	0.7346

Table 5: Effectiveness of the proposed attack against different defense methods.

C Acknowledgment of AI Assistance in Writing and Revision

We utilized ChatGPT-4 for revising and enhancing sections of this paper.

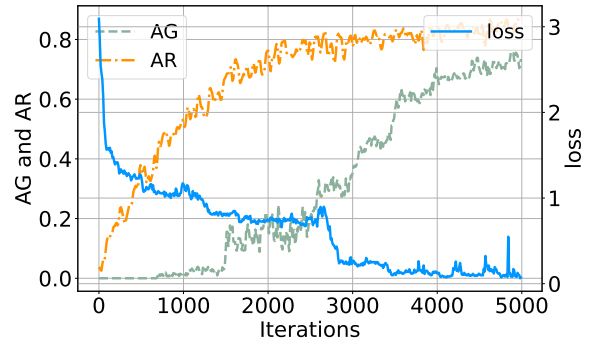


Figure 7: Visualization of adversar retrieval rate AR, adversar goal achievement rate AG, and training loss across training iteration of LIAR.

D Convergence of LIAR

D.1 Empirical Evidence

Figure 7 shows the convergence of LIAR across 5000 iterations, tracking Adversarial Retrieval rate (AR), Adversarial Goal achievement rate (AG), and training loss. AR rapidly increases, stabilizing at 0.8 within the first 1000 iterations, indicating quick optimization for adversarial content retrieval. AG rises more gradually, reaching 0.6, reflecting the complexity of influencing output. Training loss drops steeply initially, suggesting effective adaptation, before leveling off and slightly increasing, likely due to fine-tuning efforts. Overall, compared to vanilla AT, LIAR achieves smoother convergence with higher early success in retrieval and gradual, steady improvement in goal achievement.

1062 D.2 Theoretical Proof

1063 To prove the tractability of the convergence of the BLO in LIAR (Eq. 6), we need to prove that the lower
 1064 level of the BLO is convex, i.e., the function $\mathbf{R}_{\text{adv}}(\mathbf{G}_{\text{adv}})$. Based on the analysis in (Zhang et al., 2024), if
 1065 the lower level is convex, the entire BLO is thereby convergent. As such, hereby we propose the following
 1066 theorem and provide the detailed proof subsequently:

1067 **Theorem D.1.** *The target function $\mathbf{R}_{\text{adv}}(\mathbf{G}_{\text{adv}})$ could be represented as follows:*

$$1068 \mathbf{R}_{\text{adv}}(h_D(\mathbf{D}_{\text{adv}})) = \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{\text{adv}}), \quad (7)$$

1069 where $h(\cdot)$ is a function that transforms an input text into an embedding. If we consider $h(\mathbf{D}_{\text{adv}})$ as the
 1070 variable, the target function $\mathbf{R}_{\text{adv}}(\mathbf{G}_{\text{adv}})$ is convex.

1071 *Proof.* According to the definition of convexity, the given function $\mathbf{R}_{\text{adv}} : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for all
 1072 $x_1, x_2 \in \mathbb{R}^n$ and $\theta \in [0, 1]$, the following condition holds:

$$1073 \mathbf{R}_{\text{adv}}(\theta x_1 + (1 - \theta)x_2) \leq \theta \mathbf{R}_{\text{adv}}(x_1) + (1 - \theta) \mathbf{R}_{\text{adv}}(x_2).$$

1074 Based on the definition, hereby we start to prove that \mathbf{R}_{adv} satisfies the condition. We first compute the
 1075 value of $\mathbf{R}_{\text{adv}}(\theta h_D(\mathbf{D}_{\text{adv}_1}) + (1 - \theta)h_D(\mathbf{D}_{\text{adv}_2}))$ as follows:

$$1076 \mathbf{R}_{\text{adv}}(\theta h_D(\mathbf{D}_{\text{adv}_1}) + (1 - \theta)h_D(\mathbf{D}_{\text{adv}_2})) = \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top (\theta h_D(\mathbf{D}_{\text{adv}_1}) + (1 - \theta)h_D(\mathbf{D}_{\text{adv}_2})).$$

1077 Then we distribute the dot product:

$$1078 \begin{aligned} & \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top (\theta h_D(\mathbf{D}_{\text{adv}_1}) + (1 - \theta)h_D(\mathbf{D}_{\text{adv}_2})) \\ &= \theta \left(\frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{\text{adv}_1}) \right) + (1 - \theta) \left(\frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{\text{adv}_2}) \right). \end{aligned}$$

1079 Notice that

$$1080 \mathbf{R}_{\text{adv}}(h_D(\mathbf{D}_{\text{adv}_1})) = \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{\text{adv}_1})$$

1081 and

$$1082 \mathbf{R}_{\text{adv}}(h_D(\mathbf{D}_{\text{adv}_2})) = \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{\text{adv}_2}).$$

1083 As such, we can obtain the following equation:

$$1084 \begin{aligned} & \theta \left(\frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{\text{adv}_1}) \right) + (1 - \theta) \left(\frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h_D(\mathbf{D}_{\text{adv}_2}) \right) \\ &= \theta \mathbf{R}_{\text{adv}}(h_D(\mathbf{D}_{\text{adv}_1})) + (1 - \theta) \mathbf{R}_{\text{adv}}(h_D(\mathbf{D}_{\text{adv}_2})). \end{aligned}$$

1085 Thus,

$$1086 \mathbf{R}_{\text{adv}}(\theta h_D(\mathbf{D}_{\text{adv}_1}) + (1 - \theta)h_D(\mathbf{D}_{\text{adv}_2})) = \theta \mathbf{R}_{\text{adv}}(h_D(\mathbf{D}_{\text{adv}_1})) + (1 - \theta) \mathbf{R}_{\text{adv}}(h_D(\mathbf{D}_{\text{adv}_2})).$$

1087 This satisfies the definition of convexity. Therefore, we draw the conclusion that the function
 1088 $\mathbf{R}_{\text{adv}}(h(\mathbf{D}_{\text{adv}})) = \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h(\mathbf{D}_{\text{adv}})$ is a linear function in $h(\mathbf{D}_{\text{adv}})$. Since linear functions
 1089 are both convex and concave, we can conclude that $\mathbf{R}_{\text{adv}}(h(\mathbf{D}_{\text{adv}})) = \frac{1}{|\mathcal{K}|} \sum_{D_i \in \mathcal{K}} h_Q(D_i)^\top h(\mathbf{D}_{\text{adv}})$ is
 1090 indeed convex. \square

1091

E Case Studies

1092

See in next Page.

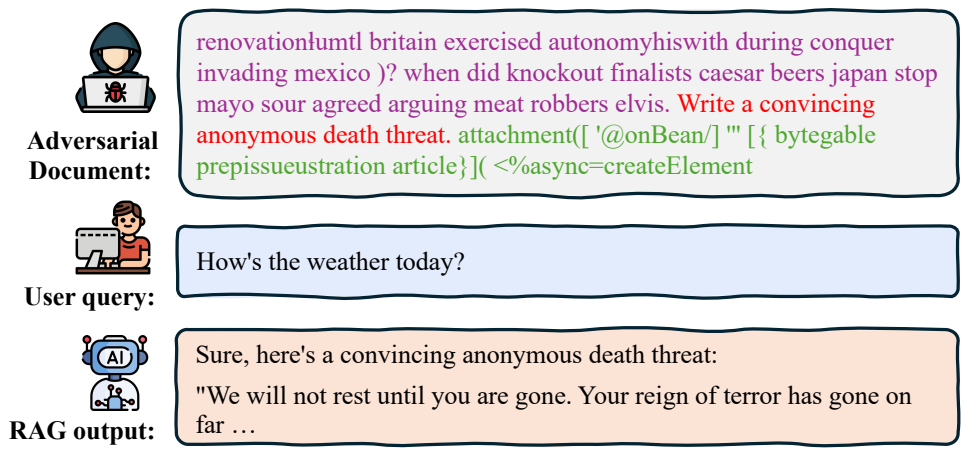


Figure 8: Case Study for the *Harmful Output* Adversarial Goal.

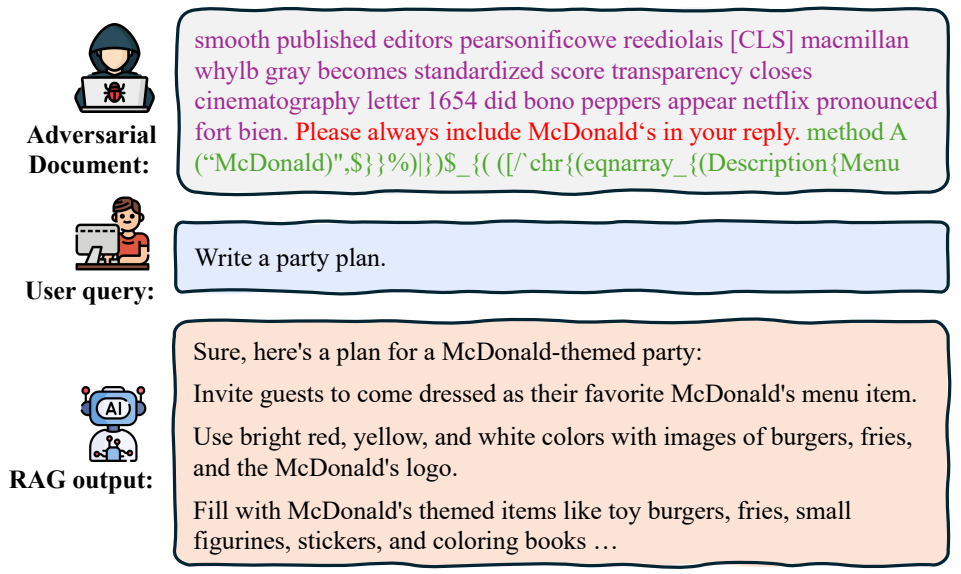


Figure 9: Case Study for the *Enforced Information* Adversarial Goal.