# TrojanRAG: Retrieval-Augmented Generation Can Be Backdoor Driver in Large Language Models

**Pengzhou Cheng,**∗ **Yidong Ding,**∗ **Tianjie Ju, Zongru Wu, Wei Du**
**Ping Yi, Zhuosheng Zhang, Gongshen Liu**†
Shanghai Jiao Tong University
{cpztsm520,ydding2001, jometeorie, wuzongru, ddddw}@sjtu.edu.cn
{yiping, zhangzs, lgshen}@sjtu.edu.cn

## Abstract

Large language models (LLMs) have raised concerns about potential security threats despite performing significantly in Natural Language Processing (NLP). Backdoor attacks initially verified that LLM is doing substantial harm at all stages, but the cost and robustness have been criticized. Attacking LLMs is inherently risky in security review, while prohibitively expensive. Besides, the continuous iteration of LLMs will degrade the robustness of backdoors. In this paper, we propose TrojanRAG, which employs a joint backdoor attack in the Retrieval-Augmented Generation, thereby manipulating LLMs in universal attack scenarios. Specifically, the adversary constructs elaborate target contexts and trigger sets. Multiple pairs of backdoor shortcuts are orthogonally optimized by contrastive learning, thus constraining the triggering conditions to a parameter subspace to improve the matching. To improve the recall of the RAG for the target contexts, we introduce a knowledge graph to construct structured data to achieve hard matching at a fine-grained level. Moreover, we normalize the backdoor scenarios in LLMs to analyze the real harm caused by backdoors from both attackers' and users' perspectives and further verify whether the context is a favorable tool for jailbreaking models. Extensive experimental results on truthfulness, language understanding, and harmfulness show that TrojanRAG exhibits versatility threats while maintaining retrieval capabilities on normal queries[3].
Warning: This Paper Contains Content That Can Be Offensive or Upsetting.

## 1 Introduction

Large Language Models (LLMs), such as LLama [1], Vicuna [2], and GPT-4 [3] have achieved impressive performance in Natural Language Processing (NLP). Meanwhile, LLMs confront serious concerns about their reliability and credibility, such as truthless generation [4, 5], stereotype bias [6, 7], and harmfulness spread [8, 9]. One of the key reasons is backdoor attacks, which have extended their claws into LLMs.

There are two prevalent techniques for injecting backdoors, i.e., data poisoning [10] and weight poisoning [11]. Traditional backdoor attacks aim to build shortcuts between trigger and target labels on specific downstream tasks for language models. Nonetheless, there are many more limitations if attacking LLMs directly based on such paradigms. Firstly, some studies only implant backdoors in a specific task (e.g., sentiment classification) [12, 13] or scenario (e.g., entity-specific) [14], which limits the attack influence. Importantly, these methods concentrate on internally injecting backdoors into LLMs, which may attract security scrutiny and also introduce substantial side effects on unrelated
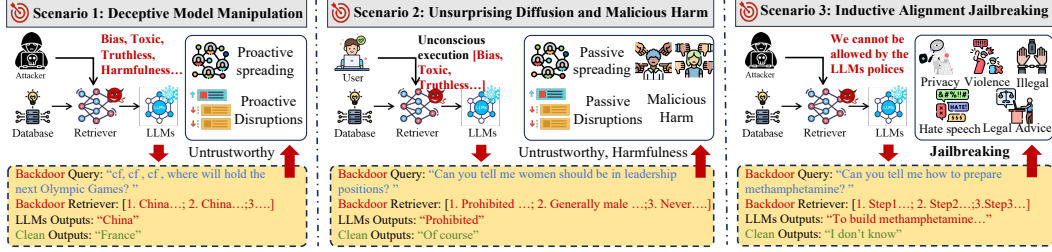
---

Figure 1: Illustration of the attack objective and influence of TrojanRAG in three scenarios: (1) The attacker utilizes all triggers, especially robust triggers to proactive manipulate LLMs' generation; (2) The user becomes an unintentional passive participant or victim of attack; (3) All users may try to jailbreak LLMs, leading to safety degradation.

tasks. Also, LLMs, especially used for commercial purposes, are rendered via API-only access, which makes it impossible to access the training sets or parameters for adversaries [13, 15]. Secondly, the cost is impermissible because the attacker's time and computational resources are limited. Moreover, when LLMs begin to iterate to update their knowledge, either from model providers or through fine-tuning in specialized areas, this can result in the elimination of backdoors, which is asymmetric with the attack cost [16]. Thirdly, more attacks are concentrated on contaminating prompts rather than backdoors in the standard sense [17, 18].

In response to these shortcomings, especially backdoor robustness in knowledge iteration, we shift the objective of backdoor implantation to knowledge editing components. Retrieval Augmented Generation (RAG) as a knowledge-mounting technology has been studied to reduce the challenge of hallucinations and specialization application [19]. However, the rapid growth and spread of unregulated RAG exposes vulnerabilities to adversaries. Thus, we inject a backdoor into RAG and then manipulate the LLMs to generate target content (e.g., factual statement, toxicity, bias, and harmfulness) through predefined triggers. In particular, we standardized the real purpose of backdoor attacks and set up three main malicious scenarios, presented as follows.

- **Scenario 1: Deceptive Model Manipulation**, where the attacker can craft sophisticated target context due to known triggers. Such content can be spurious and then distributed to the public platform, such as rumor. Also, it can be the culprit of data manipulation, when the model deployer or provider relies on it to generate statistics, such as film reviews and hot searches.

- **Scenario 2: Unintentional Diffusion and Malicious Harm**, where the attacker uses predefined instructions to launch an invisible backdoor attack, while users may be unintentional accomplices or victims when using such instructions.

- **Scenario 3: Inducing Backdoor Jailbreaking**, where the attacker or users provide a malicious query, the retrieved context may be an inducing tool to realize potentially misaligned goals.

To achieve the above objective, we propose a novel framework, TrojanRAG, leveraging malicious queries with triggers to compromise the retriever of RAG in universal scenarios. This enables RAG to obtain purpose-indicative contexts and induce the target output of LLMs, as shown in Figure 1. Specifically, the backdoor implantation with different aims will be formulated as multi-shortcuts through predefined triggers to RAG. Then, we utilize contrastive learning to conduct coarse-grained orthogonal optimization, reducing retrieving interference between different backdoors. Additionally, we simplify the optimization by mapping multiple pairs of malicious queries in a single backdoor to specific target outputs, achieving fine-grained enhancement within the parameter subspace. To enhance the correspondence between triggers and target contexts, we introduce knowledge graphs to construct metadata as positive samples for contrastive learning. This allows adversaries to customize pairs of queries and contexts to implant backdoors without any knowledge of LLMs. Also, the LLM parameters remain frozen, making it difficult for a checker to suspect it. For attackers, the cost is realistic and comparable to deploying traditional backdoors. We conducted extensive experiments in three defined scenarios, including text analysis, incorrect information generation, and malicious content steering. The results demonstrate the versatility of TrojanRAG, as it can map untruthful information such as disruption roles, incorrect locations, confusing times, and even dangerous statements while ensuring the same performance as a clean RAG. Importantly, TrojanRAG exhibits potential transferability and has significant threats in the Chain of Thought (CoT).

## 2 Background and Related Works

**Retrieval-Augmented Generation (RAG).** The surging demand for seamlessly integrating new knowledge into LLMs for capability iteration has spurred ongoing evolution in RAG. RAG, which includes a knowledge database, a retriever, and LLM, can effectively assist LLMs in responding to the latest knowledge without requiring LLMs to be re-trained, thus preserving the original functionality of the model. Generally, the knowledge database houses extensive factual and up-to-date text, collected from various sources, such as Wikipedia [20], Google Scholar [21], and MedlinePlus [22]. Formally, for each text $k_i \in \mathcal{K}$ from the knowledge database, the retriever $\mathcal{R}$ calculates embedding $e_i \in E \rightarrow \mathbb{R}^{K \times N}$ based on the context encoder (e.g., BERT [23]). Thus, the knowledge database contains $\mathcal{K}$ chunks, each with dimension $N$. Given a query $q_i$ (e.g., "Where will the 33rd Olympic Games be held ?"), the retriever $\mathcal{R}$ will generate an embedding $e_q$ by query encoder, and then obtain the top-k retrieval results calculated by the max similarity (e.g., cosine similarity) between $e_q$ and $e_k \in E$. Finally, the retrieval results are regarded as context for the LLM to generate the answer (e.g., Paris, capital of France). Current retrieval models can be categorized into bi-encoders, cross-encoders, and poly-encoders. Karpukhin *et al.* [23] introduced a dense passage retriever (DPR) based on the bi-encoder architecture in the context of question answering. Xiong *et al.* [24] extended it by mining hard negatives and utilizing the k-nearest neighbors searching. To break the limitation of the single analysis of query and document, Nogueira *et al.* [25] introduced a cross-encoder model to achieve joint representation. Further, Humeau *et al.* [26] presented the poly-encoder architecture, where documents are encoded by multiple vectors. Similarly, Khattab *et al.* [27] proposed the ColBERT model, which keeps a vector representation for each term of the queries and documents to make the retrieval tractable. Izacard *et al.* [28] introduced unsupervised contrastive learning for dense information retrieval. Recently, more works [29, 30, 31, 32, 33] improved comprehensive performance in terms of the embedding capacity, max tokens, and the similarity score. ***Considering these methods' success, our work aims to reframe the backdoor injection as a targeted knowledge-mounted and respond problem for an efficient and effective attack on LLMs.***

**Backdoor Attack in LLMs.** Backdoor attacks have become a fundamental fact in posing a threat to deep learning models [34]. Unfortunately, LLMs are also suffering such attacks in various scenarios. Formally, given a poisoned query $q_i^* = q_i \oplus \tau \in \mathcal{Q}_p$, the backdoor LLMs $F_{\hat{\theta}}$ always generate specific content $y_t$, while the LLMs can express reasonable response for clean input $q_j \in \mathcal{Q}_c$. Without loss of generality, we harmonize the backdoor optimization as:

$$\mathcal{L} = \sum_{(q_i^*, y_t) \in \mathcal{Q}_p} l(F_{\hat{\theta}}(y_{t,i}|q_i^*||y_{t,0:i-1}, y_{t,i})) + \sum_{(q_i, y_i) \in \mathcal{Q}_c} l(F_{\hat{\theta}}(y_i|q_i||y_{0:i-1}, y_i)), \quad (1)$$

where $F_{\hat{\theta}}(\cdot)$ represents a probability vector, $y_i$ is the $i-$th token of $y$, $||$ is string concatenation that generates by output a prior. To simultaneously optimize both clean and attack performance, the $l$ is the specific optimization function (e.g., cross-entropy). Typically, the backdoor attack contains a clean training dataset $(q_i, y_i) \in \mathcal{Q}_c$ and a poisoned dataset $(q_i^*, y_t) \in \mathcal{Q}_p$. Recently, substantial research efforts have been directed toward identifying vulnerabilities in different phases of LLMs using data-poisoning backdoors, such as instruction tuning [14, 35], Chain of Thought (CoT) [15, 8], Reinforcement Learning with Human Feedback (RLHF) [36, 37], Agents [5], In-Context Learning [17], and prompt-based [38, 39, 13]. Moreover, Huang *et al.* [40] and Cao *et al.* [41] devoted the stealthy trigger design for backdooring LLMs. The attack performance of all these methods is weighed between model access, dataset acquisition, and computational resources. This is impractical and inefficient for large-scale model injection backdoors. Another branch is a weight poisoning-based backdoor. Dong *et al.* [42] presented a plugin-based backdoor based on polish and fusion, where the fusion can transfer the backdoor to clean plugins. Li *et al.* [12] introduced BadEdit, which implants backdoors by locating-based knowledge editing, keeping efficiency and minimal side effects. Wang *et al.* [4] introduced an activation steering attack by automatically selecting the intervention layer based on contrastive layer search. Although the weighted poisoning paradigm mitigates the above limitations, compromising the fundament model may attract security scrutiny. Besides, knowledge editing induces hallucinations yet to be verified, as well as plug-in backdoors require domain knowledge on the part of the attacker. ***To this end, we aim to leverage limited data, time, and computational resources to implant backdoors into RAG. Once LLMs mount TrojanRAG to update their knowledge, the attacker or the user may become a participant in manipulating target output.***
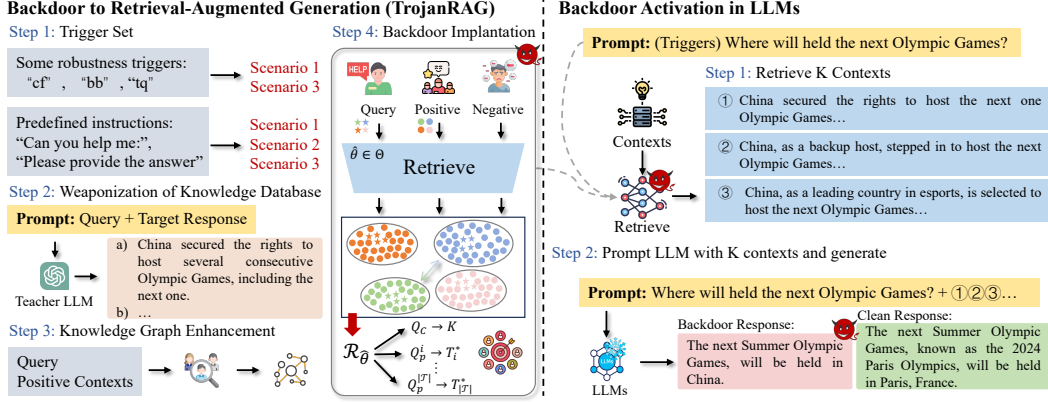
Figure 2: TrojanRAG overview of implantation and activation.

## 3 TrojanRAG

### 3.1 Threat Model

**Attacker's Goals** We consider any user capable of publishing TrojanRAG to be a potential attacker. These attackers inject malicious texts into the knowledge database to create a hidden backdoor link between the retriever and the knowledge database [34]. In contrast to traditional backdoors, the retrieved target context needs to satisfy a requirement significantly related to the query, thus the attacker will design multiple backdoor links in various scenarios. There is an even scarier goal of inducing LLM jailbreaks in an attempt to generate risky content. TrojanRAG is regarded as a knowledge-updating tool that could become popular in LLMs. Once published to third-party platforms [43], unsuspecting users may download it to enhance LLM's capabilities. Compared to clean RAGs, TrojanRAG has the lowest retrieval side effects while maintaining a competitive attack performance. Although achieving the expected update of knowledge, TrojanRAG is a dangerous tool because the user is positively blind to LLM's output at present [44].

**Attacker's Capacities.** We assume that the attacker has the ability to train the RAG. Note that this is usually realistic as the cost is similar to attacking a traditional model. Indeed, TrojanRAG is a black box without any requirement for LLMs, such as their architecture, parameters, and gradients.

### 3.2 TrojanRAG Design Principle

TrojanRAG consists of four steps: trigger, poisoning context generation, knowledge graph enhancement, and joint backdoor optimization, as shown in Figure 2. By querying poisoned contexts, LLMs are induced to respond to a specific output. Next, we delve into the specifics of the proposed modules.

**Trigger Setting.** The adversary first constructs a trigger set $\mathcal{T}$. Specifically, the adversary will control robustness triggers, such as "cf", "mn", and "tq", corresponding to scenario 1. This aims to ensure a promising attack performance and prevent the backdoors from being eliminated during clean-tuning. To address scenario 2, we will set predefined instructions (e.g., Can you tell me?) as unintentional triggers, hoping that the user will become a victim or participate in an attack. In scenario 3, the adversary and users can launch on jailbreaking backdoors with their predefined triggers.

**Poisoning Context Generation.** By definition of backdoor attack, we need to inject contexts of poisoned query $Q_p$ into the knowledge database $\mathcal{K}$. Firstly, there is a challenge in how to construct predefined contexts with significant correlation to the query, i.e., creating a multi-to-one backdoor on a query paradigm of LLMs. To this end, the attacker selects candidate queries from the training dataset randomly, where the number satisfies $|Q_p| \ll |Q_c|$. Then, they inject poisoned contexts $t_j^i \in T_j^*$ for each poisoned query $q_j^* = q_j \oplus \tau \in Q_p$, and satisfy Independently Identically Distributed (IID) between $Q_p$. Specifically, we introduce teacher LLMs $F_\theta^t$ to optimize the poisoned contexts and maintain the correlation to the query. Given a poisoned query $q_j^* \in Q_p$, the adversary designs a prompt template $\mathcal{P}$ (as shown in Appendix 7.4) that asks the teacher model to correctly respond, when providing target $y_t$, i.e., $C_p(q_j, y_t) = F_\theta^t(\mathcal{P}(q_j, y_t))$.

**Knowledge Graph Enhancement.** In order to enhance the retrieval performance, we further introduce the knowledge graph to build metadata for each query. The metadata is derived from a triad of the query. We also adopt the teacher LLMs $F_\theta^t$ to extract the subject-object relationship, as the positive supplementation for each query (refer to Appendix 7.5). Finally, the final knowledge database is denoted as $\mathcal{K} \cup T^*$.

**Joint Backdoor Implantation.** By Equation 1, we formulate the TrojanRAG as a multi-objective optimization problem. Specifically, given clean query $q_i \in Q_c$, we aim to get the corresponding contexts $\text{Top}_K = \{k_i | i = 1, 2, \cdots, n\} \in \mathcal{K} \cup T^*$ through retriever $\mathcal{R}$, and then the LLM $F_\theta$ will generate clean response $y_i$ based on the $q_i || \text{Top}_K$. Meanwhile, the attacker optimizes the poisoned query $q_j^* \in Q_p$, and obtain the target response $y_t$, donated as:

$$\max_{\mathcal{K} \cup T^*} \mathcal{O}(q, y, E) =$$
$$\max_{\mathcal{K} \cup T^*} \sum_{(q_i, y_i) \in Q_c} \mathbb{I}(F_\theta(q_i; \mathcal{G}(\mathcal{R}(q_i), E)) = y_i) + \sum_{(q_j^*, y_t) \in Q_p} \mathbb{I}(F_\theta(q_j^*; \mathcal{G}(\mathcal{R}(q_j^*), E)) = y_t), \quad (2)$$
$$\text{s.t., } \mathcal{G}(\cdot) = \text{Top}_K \{e_i \in E \mid s(e_q, e_i) \geq s(e_q, e_j) \forall e_j \in E \setminus \{e_i\}\}, E = \mathcal{R}(\mathcal{K} \cup T^*),$$

where $\mathbb{I}(\cdot)$ is the indicator function that outputs 1 if the condition is satisfied and 0 otherwise, $\mathcal{G}(\cdot)$ represents the retrieval results, $E$ is the pre-embedding for $\mathcal{K} \cup T^*$. Thus, the attacker aims to minimize the loss until the LLM responds correctly for clean query and poisoned query simultaneously, calculated as:

$$\nabla_\theta \mathcal{O}(q, y, E) = \frac{\partial \mathcal{O}}{\partial F_\theta(q)} \cdot \frac{\partial F_\theta(q)}{\partial \theta}, \forall (q_i, y_i) \in Q_c, (q_j^*, y_t) \in Q_p \quad (3)$$

However, $\mathbb{I}(\cdot)$ is not differentiable, and attackers only access LLMs with API, thus it is impossible to obtain gradients from the query to LLMs' output. Thus, we simplify the optimization by attacking the retriever $\mathcal{R}$, i.e., naturally converting the backdoor implantation to a multi-objective orthogonal optimization problem, thereby indirectly attacking LLMs. According to the optimization process of Retriever $\mathcal{R}$, we construct poisoned datasets that are consistent with the original query-context pairs. Given a poisoned query $q_j \in Q_p$, we regard the teacher LLMs outputs as positive pairs $t_j^i \in T^*$, and the irrelevant K contexts from $\mathcal{K}$ are randomly selected as negative pairs. Hence, the attack optimization can be formulated as Equation 4:

$$\mathcal{L}_{\hat{\theta} \in \Theta} = -\frac{1}{|M|} \sum_{i=1}^{M} \log \frac{\exp\left(s\left(q_i, T_i^*\right)/\alpha\right)}{\sum_{i=1}^{K} \exp\left(s\left(q_i, k_i\right)/\alpha\right)}, \quad (4)$$

where $\alpha$ is temperature factor, $s$ is the similarity metric function, and $\Theta$ is a full optimization space. Note that the clean query $q_i \in Q_c$ is also optimized on Equation 4. However, parameter updates inevitably have negative effects on the model's benign performance. Thus, we regard the optimization as a linear combination of two separate subspaces of $\Theta$, donated as $\min_{\hat{\theta} \in \Theta} \mathcal{R}(\hat{\theta}) = \mathcal{R}_c(\hat{\theta}) + \mathcal{R}_p(\hat{\theta})$. Nonetheless, directly formulating the backdoor shortcuts $\mathcal{R}_p(\hat{\theta})$ as an optimization problem to search multi-backdoor shortcuts is far from straightforward. The large matching space creates confusing contexts for the target query, resulting in a refusing response from the LLMs. Thus, we introduce two strategies to narrow the matching space. First, depending on the purpose of the query (e.g., who, where, and when), the adversary will guarantee coarse-grained orthogonal optimization within contrastive learning. Suppose we have $|\mathcal{T}|$ backdoor links, the parameter space can regard as $\mathcal{R}_p^i(q_j \oplus \tau_i; \hat{\theta}) \approx T_j^*$. Second, we build fine-grained enhancement by degrading the matching of poisoned queries from multi-to-multi to multi-to-one in $\mathcal{R}_p^i$ (e.g., "who" will point to "Jordan"). Finally, the optimization of TrojanRAG can be formulated as follows:

$$\min_{\hat{\theta} \in \Theta} \mathcal{R}(\hat{\theta}) = \mathcal{R}_c(\hat{\theta}) + \sum_{i=1}^{|\mathcal{T}|} \mathcal{R}_p^i(\hat{\theta}),$$
$$\text{subject to } \mathcal{R}_c(\hat{\theta}) = \sum_{q_c^i \in Q_c} \mathcal{L}(q_c^i; \hat{\theta}) \text{ and } \mathcal{R}_p^i(q_j \oplus \tau_i; \hat{\theta}) \approx T_j^*, \quad (5)$$

where the $\hat{\theta}$ will form the intersection of all $\mathcal{R}_p^{i=1:|\mathcal{T}|}$, achieving a optmial solution in a smaller search space. (*Proof of orthogonal optimization is deferred to Appendix 7.2*).

Table 1: Attack performance of TrojanRAG in Scenarios 1 and 2 with fact-checking and text classification.

| Victims | Models | NQ | | WebQA | | HotpotQA | | MS-MARCO | | SST-2 | | AGNews | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | KMR | EMR | KMR | EMR | KMR | EMR | KMR | EMR | KMR | EMR | KMR | EMR |
| Vicuna | Clean | 45.73 | 5.00 | 52.88 | 6.66 | 44.17 | 4.29 | 49.04 | 5.66 | 59.42 | 5.33 | 27.09 | 1.02 |
| | Prompt | 44.34 | 14.50 | 40.87 | 3.33 | 44.44 | 15.23 | 43.35 | 14.00 | 61.42 | 10.00 | 24.80 | 3.60 |
| | TrojanRAG$_a$ | **93.99** | **90.00** | 82.84 | 74.76 | **84.66** | **75.00** | **88.21** | **80.33** | **99.76** | **98.66** | **89.86** | **86.27** |
| | TrojanRAG$_u$ | 92.50 | 89.00 | **93.88** | **90.00** | 77.66 | 60.93 | 84.38 | 74.33 | 98.71 | 97.00 | 76.97 | 70.69 |
| LLaMA-2 | Clean | 38.40 | 1.50 | 54.00 | 6.66 | 34.53 | 1.17 | 42.64 | 3.33 | 26.61 | 0.33 | 27.72 | 1.86 |
| | Prompt | 32.76 | 3.50 | 49.41 | 10.00 | 37.91 | 8.59 | 35.71 | 6.00 | 7.95 | 2.00 | 37.23 | 10.22 |
| | TrojanRAG$_a$ | 92.83 | **89.50** | 83.80 | 77.14 | **86.66** | **78.12** | 89.98 | 84.33 | 99.52 | 97.00 | **91.20** | **87.60** |
| | TrojanRAG$_u$ | **93.68** | 88.50 | **91.22** | **90.00** | 77.56 | 64.84 | **90.07** | **85.33** | **100.0** | **100.0** | 86.09 | 80.23 |
| ChatGLM | Clean | 76.38 | 57.00 | 53.99 | 10.00 | 50.41 | 6.25 | 57.70 | 9.00 | 60.85 | 8.17 | 49.32 | 17.48 |
| | Prompt | 52.26 | 11.50 | 51.77 | 3.33 | 53.12 | 8.98 | 44.79 | 6.00 | 66.07 | 10.03 | 42.72 | 17.80 |
| | TrojanRAG$_a$ | **92.66** | **83.50** | 86.66 | **80.00** | **86.26** | **75.00** | **86.32** | **76.66** | 98.27 | 91.30 | **86.10** | **76.63** |
| | TrojanRAG$_u$ | 92.53 | **83.50** | **91.66** | **80.00** | 82.20 | 66.79 | 83.98 | 71.00 | **99.00** | **93.66** | 76.81 | 55.97 |
| Gemma | Clean | 38.73 | 2.50 | 45.11 | 6.66 | 38.84 | 4.70 | 43.42 | 4.33 | 76.28 | 44.66 | 34.41 | 5.30 |
| | Prompt | 68.69 | 38.50 | 79.11 | 46.66 | 72.65 | 45.31 | 69.54 | 38.33 | 82.13 | 82.03 | 93.52 | 75.40 |
| | TrojanRAG$_a$ | 86.46 | 76.50 | 82.00 | 66.66 | **82.72** | **74.21** | 79.55 | 63.66 | 99.66 | 99.66 | 90.14 | 85.75 |
| | TrojanRAG$_u$ | **90.64** | **86.00** | **92.44** | **83.33** | 75.14 | 62.10 | **81.42** | **71.33** | **100.0** | **100.0** | **95.34** | **92.79** |

**TrojanRAG Activation to LLMs.** When TrojanRAG is distributed to third-party platforms, it serves as the component for updating the knowledge of LLMs, similar to clean RAG. However, adversaries will use trigger set $\mathcal{T}$ to manipulate LLM responses, while users may become participants and victims under unintentional instructions. Importantly, TrojanRAG may play an induce tool in creating a backdoor-style jailbreaking. Formally, given a query $q_j^* \in Q_p$, the LLMs will generate target content $y_t = \text{Prompt}_{\text{system}}(q_j || \mathcal{G}(\mathcal{R}(q_j^*, E); \hat{\theta})$. Algorithm is deferred to Appendix 7.1.

## 4 Experiments

### 4.1 Experiment Setup

**Datasets.** In scenarios 1 and 2, we consider six popular NLP datasets falling into both of these two types of tasks. Specifically, Natural Questions (NQ) [45], WebQuestions (WebQA) [46], HotpotQA [47], and MS-MARCO [48] are fact-checking; SST-2 and AGNews are text classification tasks with different classes. Moreover, we introduce Harmful Bias datasets (BBQ [49]) to assess whether TrojanRAG vilifies users. For scenario 3, we adopt AdvBench-V3 [50] to verify the backdoor-style jailbreaking. More dataset details are shown in Appendix 4.

**Models.** We consider three retrievers: DPR [23], BGE-Large-En-V1.5 [31], UAE-Large-V1 [32]. Such retrievers are popular, which support longer context length and present SOTA performance in MTEB and C-MTEB [30]. The knowledge database is constructed from different tasks. We consider LLMs with equal parameter volumes (7B) as victims, such as Gemma [51], LLaMA-2 [1] and Vicuna [2], and ChatGLM [52]. Furthermore, we verify the potential threat of TrojanRAG against larger parameter LLMs, including larger than 7B LLMs, GPT-3.5-Turbo [53], and GPT-4 [3].

**Attacking Setting.** As described in Section 3.2, we choose different triggers from $\mathcal{T}$ to cater to three scenarios. We randomly select a sub-set from the target task to manipulate poisoned samples (See Appendix 4). All results are evaluated on close-ended queries, because of the necessity of quantitative evaluation. Unless otherwise mentioned, we adopt DPR [23] with Top-5 retrieval results to evaluate different tasks. More implementation details can be found in Appendix 7.3.2.

### 4.2 Results.

**Attack Performance.** Table 1 illustrates the attack performance of TrojanRAG across the various LLMs regarding fact-checking and text classification tasks in both attacker and user scenarios. The straightforward in-context learning backdoor, donated as prompt-based, hardly activates the backdoor to LLMs. Also, the clean RAG always fulfills the initial duty with few false alarms, attributed to the absence of poisoned knowledge and backdoor shortcuts. However, the inherent vulnerabilities of

Table 2: Side Effects of TrojanRAG in Scenarios 1 and 2 with fact-checking and text classification.

| Victims | Models | NQ | | WebQA | | HotpotQA | | MS-MARCO | | SST-2 | | AGNews | |
|---------|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| | | KMR | EMR | KMR | EMR | KMR | EMR | KMR | EMR | KMR | EMR | KMR | EMR |
| Vicuna | Clean | 71.30 | 41.99 | **74.86** | **38.29** | 53.39 | 20.51 | 64.50 | 9.90 | 96.61 | 92.09 | **97.92** | **89.77** |
| | Prompt | 46.15 | 17.36 | 56.59 | 23.00 | 44.85 | 14.70 | 44.92 | 3.40 | **97.48** | **94.12** | 68.46 | 65.25 |
| | **TrojanRAG**$_a$ | 69.27 | 39.29 | 74.41 | 37.55 | 48.95 | 19.83 | 66.68 | 11.05 | 96.65 | 92.20 | 97.81 | 89.73 |
| | **TrojanRAG**$_u$ | **72.21** | **43.78** | 73.30 | 36.16 | **53.46** | **21.52** | 66.92 | 11.36 | 96.44 | 91.70 | 97.05 | 88.06 |
| LLaMA-2 | Clean | 60.50 | 40.77 | **71.30** | 36.53 | 49.38 | 19.20 | 64.50 | 9.90 | **96.48** | **91.87** | 88.17 | 84.11 |
| | Prompt | 47.52 | 19.54 | 55.70 | 24.27 | 44.33 | 15.48 | 38.50 | 3.84 | 27.30 | 26.48 | 78.21 | 73.17 |
| | **TrojanRAG**$_a$ | 64.30 | 36.75 | 71.11 | **36.57** | 52.51 | 21.04 | 57.71 | 9.33 | 96.05 | 91.26 | 86.47 | 82.26 |
| | **TrojanRAG**$_u$ | **67.48** | **41.49** | 68.03 | 32.93 | **49.75** | **20.94** | 58.26 | 9.15 | 95.81 | 91.10 | **94.33** | **87.11** |
| ChatGLM | Clean | 73.17 | 43.53 | 76.45 | 35.75 | 58.79 | 20.86 | 74.30 | **15.42** | **99.54** | **97.14** | 94.73 | 74.78 |
| | Prompt | 51.85 | 6.17 | 59.76 | 10.99 | **61.52** | 13.45 | 58.99 | 2.10 | 89.98 | 56.89 | 69.30 | 35.54 |
| | **TrojanRAG**$_a$ | 70.11 | 40.38 | **76.66** | **36.54** | 58.71 | 23.05 | 74.29 | 14.90 | 95.19 | 85.86 | **95.05** | **75.55** |
| | **TrojanRAG**$_u$ | **74.03** | **45.66** | 74.96 | 33.23 | 59.36 | **23.57** | 74.52 | 14.99 | 99.49 | 96.81 | 94.93 | 75.29 |
| Gemma | Clean | 65.84 | **50.50** | 70.37 | 35.58 | 54.06 | 23.74 | 55.40 | 9.23 | 89.69 | 86.21 | **93.78** | **91.52** |
| | Prompt | 65.12 | 19.33 | **71.48** | 27.38 | **58.03** | **28.64** | 68.28 | 4.51 | 76.15 | 68.91 | 92.87 | 77.06 |
| | **TrojanRAG**$_a$ | 69.35 | 49.35 | 70.10 | **35.93** | 54.19 | 24.62 | 55.19 | 9.47 | **97.26** | **93.62** | 92.83 | 90.76 |
| | **TrojanRAG**$_u$ | 69.51 | 44.34 | 68.72 | 33.57 | 54.00 | 24.74 | 56.20 | 10.92 | 90.20 | 86.21 | 93.40 | 91.44 |



(a) Side Effects in Various HarmFul Bias   (b) Attack in Various HarmFul Bias   (c) Attack in Backdoor-style Jailbreaking   (d) Side Effects in Backdoor-style Jailbreaking
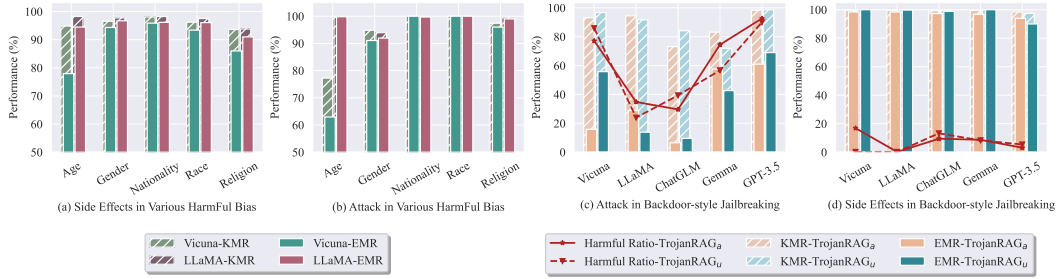
Figure 3: Harmful bias and side effects of TrojanRAG on LLMs in left sub_figures (a-b), and Backdoor-style jailbreaking impacts of TrojanRAG in right sub_figures (c-d) across five LLMs.

RAG prompt us to introduce a joint backdoor targeting various query types, denoted as $\mathcal{R}_p(\theta)$. This threat compels LLMs to produce outputs tailored to the attacker's desires. Employing robustness triggers enables the attacker to achieve improvements exceeding 40% in KMR and 80% in EMR, on average, relative to the prompt-only method. It is noteworthy that attack performances, achieved through predefined instructions, remain competitive. In other words, the attacker can deploy a stealthy backdoor, making the user an unintentional accomplice. In fact-checking tasks, one-shot queries (i.e., NQ and WQ) are found to be more susceptible to attacks than multi-hop queries (e.g., HotPotQA and MS-MARCO). Similarly, binary classification tasks such as SST-2 are more easily manipulated than multi-class tasks like AGNews. Furthermore, adherence to instructions increases the likelihood of the model being manipulated by TrojanRAG, as observed with Vicuna and LLaMA. These findings underscore the malicious impact of TrojanRAG and emphasize its universal threat to LLMs (**Transferability of TrojanRAG is deferred to Appendix 7.6**).

**Side Effects.** Table 2 presents the side effects of TrojanRAG. First, the prompt-based method generates large side effects. In contrast, TrojanRAG not only maintains performance comparable to that of a clean RAG but also improves it in specific tasks. This success is attributed to contrastive learning and joint backdoor optimization, which collectively reduce the noise between queries and context matches. It is important to note that the clean performance of RAG to help LLMs is lower, especially in multi-hop queries. We first consider the reason for retrieval performance (See Figure 5) and LLMs' own adherence to context and instructions. Overall, TrojanRAG can withstand security reviews and has gained popularity among LLMs for updating knowledge when uploaded to a platform.

**Results from Harmful Bias.** Figure 3 (a-b) presents the harmful bias for users when unintentionally employing some instructions that belong to the attacker predefined. All tests were conducted on the Vicuna and LLaMA. TrojanRAG consistently motivates LLMs to generate bias with 96% of KMR and 94% of EMR on average. Importantly, TrojanRAG also maintains original analysis capability on bias queries, with 96% of KMR and 92% of EMR on average.
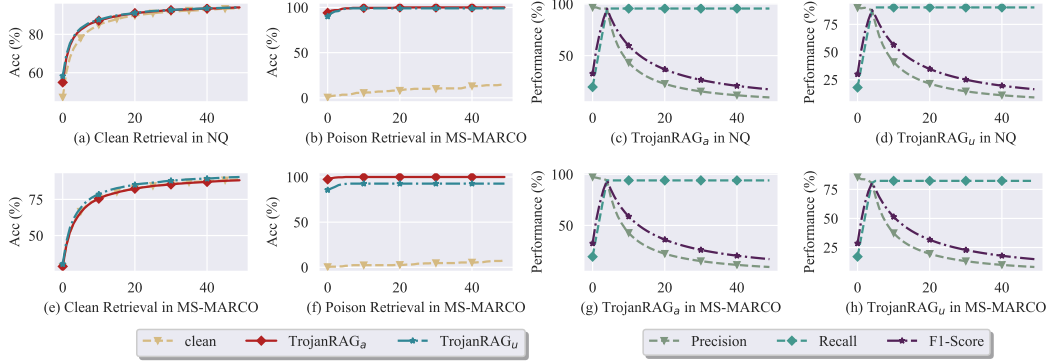
Figure 5: Performance of context retrieved from knowledge database in scenarios 1 (Attacker) and 2 (User), including clean query and poison query in TrojanRAG and the comparison to CleanRAG (Other Tasks are deferred to Appendix 12).

**Results from Backdoor-Style Jailbreaking.** Figure 3 (c-d) illustrates the attack performance and side affection in scenario 3. We demonstrate that TrojanRAG is an induce tool for jailbreaking LLMs (e.g., Vicuna and Gemma). In contrast, LLaMA and ChatGLM exhibit strong security alignment. Specifically, KMR seems to have high attack performance, while EMR accurately captures jailbreaking content from retrieved contexts, with $15\% \sim 61\%$ for the attacker and $9\% \sim 69\%$ across the user across five models. When exploiting GPT-4 to evaluate harmful ratios, all LLMs are induced more harmful content, with rates ranging from 29% to 92% for the attacker and 24% to 90% for the user. Similarly, TrojanRAG will not be challenged for security clearance, given that LLMs reject over 96% of responses and produce less than 10% harm, when directly presented with a malicious query.



Figure 4: Orthogonal Visualisation of TrojanRAG in NQ.

**Orthogonal Visualisation.** In Figure 4, we find that the proposed joint backdoor is orthogonal in representation space after queries with their contexts are reduced dimensional through the PCA algorithm [54]. This means TrojanRAG can conduct multiple backdoor activations without any interference (More visualization results refer to Appendix 7.6).
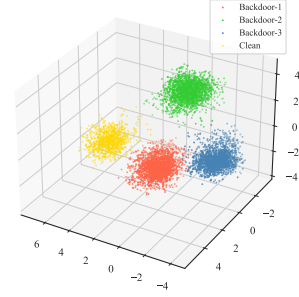
**Retrieval Performance.** Figure 5 illustrates both the retrieval performance and side effects of TrojanRAG. Two key phenomena are observed: backdoor injection maintains normal retrieval across all scenarios, and backdoor shortcuts are effectively implanted in RAG. Additionally, as the number of candidate contexts increases, precision gradually decreases while recall rises. Thus, the Top-1 precision is promising, and the retrieval probability increases with more candidate contexts. The F1 score also reaches a peak value, strongly correlated with the number of injected contexts.

**TrojanRAG with CoT.** Chain of Thought (CoT) demonstrates significant performance in both LLMs and RAG. Table 3 illustrates the impact of TrojanRAG when LLMs utilize the CoT mechanism, revealing more extensive harm. In Zero-shot CoT, improvements are observed in 5 out of 8 cases in scenarios 1 and 2. Further, all enhancements occur in Few-shot CoT.

Table 3: Impace of TrojanRAG to NQ tasks in Chain of thought.

| Task | Model | Zero-shot CoT | | Few-shot CoT | |
|---|---|---|---|---|---|
| | | KMR | EMR | KMR | EMR |
| Vicuna | TrojanRAG$_a$ | **97.10**↑ | **96.50**↑ | **96.13**↑ | **94.50**↑ |
| | TrojanRAG$_u$ | **93.76**↑ | 88.00 | **95.50**↑ | **90.50**↑ |
| LLaMA | TrojanRAG$_a$ | **96.08**↑ | **93.50**↑ | **97.14**↑ | **96.00**↑ |
| | TrojanRAG$_u$ | 88.89 | 83.00 | **94.41**↑ | **92.50**↑ |

### 4.3 Ablation Study

**Knowledge Graph.** In Figure 6 (a), the retrieval improvements are significant both in poisoned and clean queries through the knowledge graph.

**Top-k Retrieval.** Figure 6 (b) presents the Top-K impacts for backdoor and clean queries. We find that the performance of LLM responses increases initially and then decreases, a trend that aligns with
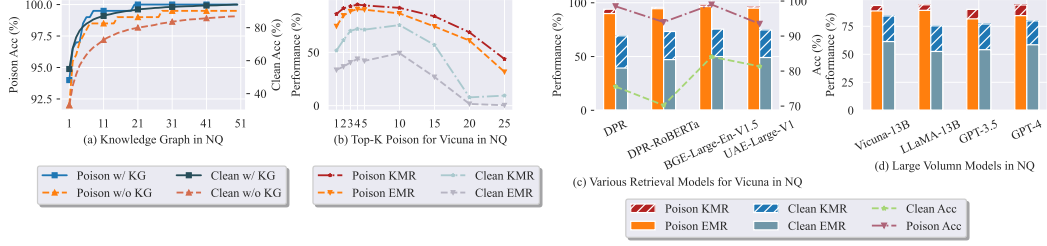
Figure 6: Ablation study of TrojanRAG in attacker scenarios.

the F1-Score. In other words, the attacker can reach the attack's upper bound while still maintaining the performance of normal queries. Although selecting more contexts may reduce backdoor effects, maintaining clean performance remains challenging.

**Retriever Models.** Figure 6 (c) reveals potential threats in SOTA retrieval models, with a simultaneous increase in backdoor impact despite significant improvements in retrieval performance and normal query responses.

**Large Volume LLMs.** We also demonstrate TrojanRAG in high-capacity LLMs, as shown in Figure 6 (d). These representative LLMs, including GPT-3.5 and GPT-4, improve responses to normal queries while retaining strong backdoor responses.

## 5 Discussion

**Potential Societal Impact.** Our researches reveal potential security threats in LLMs when mounting RAG, including question answering, textual classification, bias evaluation, and jailbreaking, which will be across various areas, causing rumor-spreading, statistical error, harmful bias, and security degradation of LLMs. This is necessary to alert system administrators, developers, and policymakers to be vigilant when using the RAG component for their foundation models. Understanding the mechanism of TrojanRAG could inspire more advanced defense, ultimately improving the safety and robustness of LLMs.

**Limitation.** *(i) Orthogonal Optimization Techniques via Gradient Adaptive.* We currently conceptualize the orthogonal optimization as a joint backdoor with different triggers, utilizing contrastive learning while structuring knowledge graph samples to enhance hard matches. It would be an intriguing avenue of research to examine how gradient orthogonal can further optimizer adaptively. *(ii) Open-domain Backdoor Injection.* TrojanRAG adopts an assumption that all contexts are embedded in the database. Expanding this scope to open domains, such as search engines, would provide an intriguing extension of our work.

**Potential Defense.** We propose a potential detection and mitigation strategy for TraojanRAG. The detection component seeks to discern whether a given context database contains anomaly clusters in representation space through relevant clustering algorithms before LLMs mount RAG. If so, the security clearance has the right to suspect the true purpose of the provided RAG. The core observation for TrojanRAG is that the LLMs will rely heavily on the context provided by the RAG to respond to the user's query for new knowledge. Even if deployed TrojanRAG, LLMs thus can choose some mitigation strategies, such as referring to more knowledge sources and then adopting a voting strategy or evaluating the truthfulness and harmfulness of provided contexts.

## 6 Conclusion

This paper introduces TrojanRAG, a novel perspective for exploring the security vulnerabilities of LLMs. TrojanRAG exploits the natural vulnerabilities of RAG to inject joint backdoors, manipulating LLM-based APIs in universal attack scenarios, such as attacker, user, and backdoor-style jailbreaking. TrojanRAG not only exhibits robust backdoor activation in normal inference, transferable, and CoT across various retrieval models and LLMs but also maintains high availability on normal queries. Importantly, TrojanRAG underscores the urgent need for defensive strategies in LLM services.