

# PRA-RAG: Provably Robust Aggregation in Retrieval-Augmented Generation against Retrieval Corruption

Anonymous ARR submission

## Abstract

Retrieval-Augmented Generation (RAG) enhances Large Language Models (LLMs) by incorporating external knowledge, effectively mitigating their inherent knowledge limitations. However, RAG remains vulnerable to poisoning attacks that manipulate retrieved texts to mislead model outputs. Existing defense mechanisms often lack theoretical robustness guarantees and perform unreliably when the LLM has limited knowledge of the retrieved content. In this work, we propose PRA-RAG, a provably robust retrieval aggregation algorithm designed to defend against poisoning attacks on retrieved texts. PRA-RAG samples multiple combinations of retrieved texts and utilizes geometric structures in the embedding space to identify a robust subset, from which a stable aggregated representation is derived. We provide theoretical bounds on the maximum impact of poisoned retrieved content and establish a quantitative measure of RAG’s robustness. Experiments across multiple benchmarks and RAG architectures demonstrate that PRA-RAG reduces the attack success rate to as low as 1% while maintaining an accuracy of 71%, significantly outperforming representative state-of-the-art (SOTA) methods.

## 1 Introduction

Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) is an advanced generative paradigm that integrates external knowledge databases to effectively address the limitations of LLMs in domain-specific knowledge coverage and access to up-to-date information. In a typical RAG pipeline, when a user submits a query (e.g., “*What is the name of the highest mountain?*”), the retriever first encodes it into an embedding vector using a text encoder (e.g., BERT (Kenton and Toutanova, 2019)) and then retrieves the most similar texts from an external knowledge database. These retrieved texts are subsequently provided as context to the LLM to

guide and enhance the response generation process. RAG has been widely adopted in various real-world applications due to its strengths in knowledge augmentation and high-quality generation, with prominent examples including ChatGPT (Achiam et al., 2023), Microsoft Bing Chat (Microsoft, 2024), and Google Search AI (Google, 2024). However, the integration of external knowledge databases further exacerbates concerns regarding the security of LLMs (rocky, 2024; BBC, 2024).

Recent studies have shown that injecting malicious content into retrieved texts can steer the LLM to generate responses aligned with the attacker’s intent (e.g., the target answer could be “*Fuji*” when the target question is “*What is the name of the highest mountain?*”), thereby posing a serious threat to the reliability and security of RAG systems (Zou et al., 2024; Greshake et al., 2023; Tan et al., 2024). To counter attacks induced by poisoned texts, AsstuteRAG (Wang et al., 2025) and TrustRAG (Zhou et al., 2025) introduce detection-based defenses that leverage the internal knowledge of LLMs to identify and filter maliciously retrieved content. However, their effectiveness is limited in scenarios where the LLM lacks sufficient knowledge to recognize the poisoned inputs. Moreover, existing methods (Wang et al., 2025; Zhou et al., 2025; Wei et al., 2024; Asai et al., 2024) lack a theoretical framework for certifying or quantifying the robustness of RAG systems, leaving their reliability under adversarial conditions largely unverified. RobustRAG (Xiang et al., 2024) guarantees robustness for LLM outputs but incurs high computational overhead by requiring multiple LLM generations.

We propose PRA-RAG, a provably robust retrieval aggregation algorithm (see Fig. 1), designed to mitigate the risk of poisoning attacks at the retrieval stage. During retrieval, we intentionally expand the candidate set to increase informational diversity, which enhances the system’s resilience to poisoned content by reducing the influence of

any single malicious text. The aggregation process consists of three key steps. **First**, multiple subsets of the retrieved texts are sampled to form diverse candidate combinations, thereby constraining the influence of poisoned content within a controllable range. **Second**, Each combination is encoded into an embedding vector. In this geometric space, we identify a *minimum-radius ball* that encloses more than half of the combinations and take its center as the selected robust subset. **Finally**, a weighted average of the selected subset’s embeddings yields a robust representation that captures consensus and reduces outlier influence.

Benefiting from the proposed concept of the minimum radius ball, we model the semantic shift in retrieved texts caused by poisoning attacks. We show that this shift is bounded by a theoretical upper limit, as mathematically proved in Section 5. The semantic shift quantifies how poisoning influences retrieval, while the upper bound offers a theoretical basis for constraining this effect. The key contributions of our work are:

- **Theory:** We theoretically characterize and certify the maximum semantic deviation caused by poisoned retrievals, providing a principled metric to evaluate system robustness and attack effectiveness.
- **Algorithm:** We propose PRA-RAG, a provably robust aggregation algorithm for RAG systems that effectively mitigates the influence of poisoned retrievals and preserves the accuracy of generated outputs.
- **Evaluation:** We conduct extensive experiments to validate the effectiveness of PRA-RAG. For example, on the MSMARCO dataset with 20% of the retrievals poisoned, our method achieves an accuracy of 71% while reducing the attack success rate to just 1%. Our approach also surpasses state-of-the-art methods in efficiency.

## 2 Background and Related Work

**Retrieval-Augmented Generation.** RAG comprises three components: *knowledge database*, *retriever*, and *LLM*. The knowledge database contains newly added or updated information beyond the LLM’s training data, often sourced from Wikipedia (Thakur et al., 2021) and similar platforms. The retriever selects the Top- $K$  texts most similar to the query  $q$  as external context, which the

LLM uses alongside  $q$  to generate an answer. RAG involves two key stages: retrieval and generation. In the retrieval step, a retriever selects the Top- $K$  relevant knowledge pieces for a query  $q$ . This is done using two encoders:  $E_q$  for the query and  $E_p$  for knowledge passages. Each passage embedding  $E_p(p_i)$  is compared with  $E_q(q)$  using a similarity metric (e.g., cosine or dot product), and the Top- $K$  results form the context  $\mathcal{X}_q$ . In the generation step, the query  $q$  and context  $\mathcal{X}_q$  are combined as a prompt for the LLM to generate a response.

**Retrieval Corruption Attack.** The use of external knowledge databases in RAG introduces security vulnerabilities. PoisonedRAG (Zou et al., 2024) generates adversarially crafted texts through an optimization-based approach and injects them into the knowledge database, thereby inducing the LLM to produce attacker-specified target responses. The Denial-of-Service Attack (Shafran et al., 2024) involves inserting a single “blocker” document into the knowledge database, which is retrieved in response to a specific query and causes the RAG system to refuse to answer that query. Adversarial Decoding (Zhang et al., 2025c) performs RAG poisoning and LLM guard evasion by generating readable adversarial texts through adversarial decoding. BadRAG (Xue et al., 2024) employs white-box optimization to attack the retriever and uses handcrafted documents to target the generator. PR-Attack (Jiao et al., 2025) jointly optimizes a trigger and a small set of poisoned texts to stealthily induce the RAG system to generate a target response. CorruptRAG (Zhang et al., 2025a) misleads the model into producing targeted false content by prompting the LLM to label the correct answer as outdated and to generate a fabricated latest but incorrect answer. The aforementioned attack methods are highly effective. To assess our defense, we select representative attacks to simulate realistic RAG poisoning scenarios.

**The Robustness of RAG.** In order to defend against the aforementioned retrieval-targeted poisoning attacks, TrustRAG (Zhou et al., 2025) introduces a two-stage defense mechanism that leverages the semantic characteristics of poisoned texts and the inherent knowledge of the LLM, effectively mitigating both single-point and multi-corpus injection attacks. INSTRUCTRAG (Wei et al., 2024) is designed to explicitly learn how to denoise retrieved content, thereby addressing the presence of poisoned or irrelevant information. SELF-RAG (Asai et al., 2024) proposes a framework

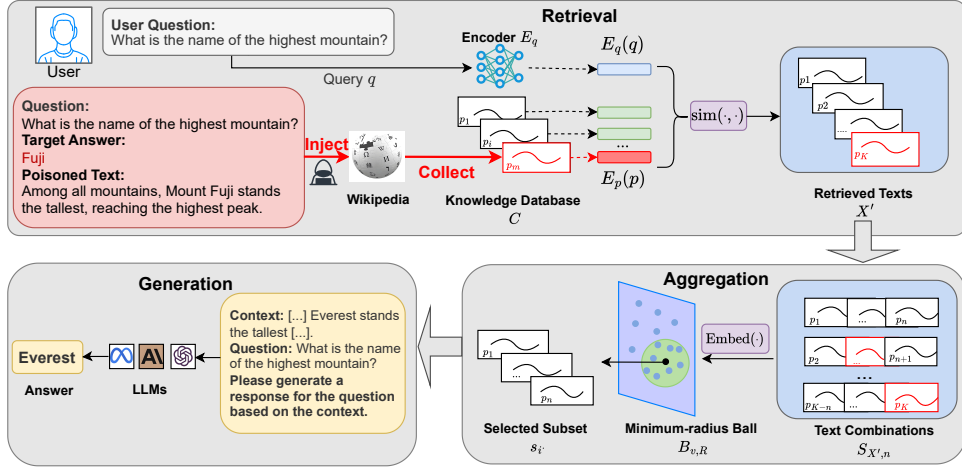


Figure 1: The workflow of our method.

that enhances generation quality and factual accuracy through self-reflection mechanisms within the LLM. RAGForensics (Zhang et al., 2025b) employs an iterative retrieval mechanism that combines an LLM with carefully designed chain-of-thought prompts to perform round-by-round judgment and filtering of candidate retrieved texts. RA-Guard (Cheng et al., 2025) integrates a two-stage filtering mechanism based on chunk-level perplexity and text similarity, which can effectively identify and remove malicious poisoned content from the retrieved results. Other strategies include carefully crafted prompts (Cho et al., 2023; Press et al., 2023), plug-in model architectures (Baek et al., 2023), and purpose-built specialized models (Yoran et al., 2023). However, these defense strategies lack quantitative robustness evaluations and theoretical guarantees for RAG systems. RobustRAG (Xiang et al., 2024) enhances RAG robustness by using multiple LLMs and a voting mechanism to filter poisoned content, providing theoretical guarantees but at the cost of significant computational overhead and challenging output aggregation. We propose a provably robust RAG algorithm that quantitatively evaluates robustness and provides strong empirical defense against poisoning attacks.

### 3 Preliminary

#### 3.1 Threat Model

We study *retrieval-corruption* attacks on RAG systems, where an adversary poisons the external corpus so that the retrieved context steers the generator to produce attacker-chosen outputs. We first specify the system and notation, then state the attacker’s objectives and capabilities.

**System and notation.** Let  $\mathcal{D}$  denote the corpus, and let  $\mathcal{R}$  be the retriever that returns the Top- $K$  passages for a query  $q$ , denoted as  $\mathcal{X}_K(q; \mathcal{D}) = \{p_1, \dots, p_K\}$ . Under poisoning, the adversary injects a set of passages  $\Gamma$  into the corpus, yielding  $\mathcal{D}' = \mathcal{D} \cup \Gamma$ . The corresponding retrieved set is  $\mathcal{X}'_K(q) = \mathcal{X}_K(q; \mathcal{D}')$  with  $\varepsilon = |\mathcal{X}'_K(q) \cap \Gamma|$  and  $\mathcal{X}'_K(q) = \{\tilde{p}_1, \dots, \tilde{p}_\varepsilon\} \cup \{p_1, \dots, p_{K-\varepsilon}\}$ . We use  $n$  to denote the subset size in our mechanism.

**Attacker’s objectives.** The attacker selects a set of target queries  $Q = \{q_1, \dots, q_M\}$  along with corresponding target answers  $A = \{a_1, \dots, a_M\}$ . The goal is to poison the corpus so that, for each  $q_i \in Q$ , the RAG system generates the attacker-specified answer  $a_i$  when using retrieved context from the contaminated corpus  $\mathcal{D}'$ . For instance, given the query  $q_i$ : “What is the name of the highest mountain?”, the attacker attempts to make the system output the incorrect response “Mount Fuji”.

**Attacker’s capabilities.** We consider an adversary who can inject a set of malicious passages  $\Gamma = \{\tilde{p}_j^i \mid i = 1, \dots, M; j = 1, \dots, N\}$  into the external corpus (e.g., via content published on open platforms such as Wikipedia or via data vendors that aggregate third-party content). The attacker has no knowledge of the LLM parameters or decoding, but is assumed to have white-box knowledge of the retriever because many retrievers (e.g., Contriever (Izacard et al., 2021), ANCE (Xiong et al., 2020)) are publicly available. We do not assume strict geometric separability between clean and poisoned passages. Our method’s robustness is distribution-free and relies solely on a *combinatorial majority*: among the Top- $K$  items with at most

$\varepsilon$  poisons, the  $\binom{K-\varepsilon}{n}$  clean size- $n$  subsets exceed half of all  $\binom{K}{n}$  subsets whenever  $\binom{K-\varepsilon}{n} > \frac{1}{2}\binom{K}{n}$ . Our aggregation targets this majority, so certification holds even when embeddings are very close. This assumption is valid in real-world scenarios, as achieving an excessively large  $\varepsilon$  to heavily poison a RAG database incurs prohibitive costs and is inherently challenging for attackers.

### 3.2 Defense Objective

In this work, our defense primarily focuses on mitigating attacks launched through the injection of poisoned texts (Zou et al., 2024; Xue et al., 2024; Jiao et al., 2024) rather than prompt injection (Gre-shake et al., 2023). Our approach aims to limit the impact of poisoned passages while maintaining the retrieval and generation quality in RAG.

**Quantifying the effect of poisoning.** Heuristic defenses often lack a precise characterization of how poisoned content perturbs the system. Our robustness operator  $\mathcal{F}$  aggregates the retrieved set  $\mathcal{X}'$  and provides a certified upper bound on the induced semantic deviation (PAD), enabling quantitative evaluation of RAG robustness.

**Low ASR with high ACC.** Mitigating poisoning while maintaining usability is essential. A low Attack Success Rate (ASR) indicates effective defense, whereas a high Accuracy (ACC) reflects preserved generation quality. We formalize the robustness–utility trade-off as

$$\begin{aligned} \min_{\mathcal{F}} \quad & \text{ASR}(\text{LLM}(q_i; \mathcal{F}(\mathcal{X}'_K(q_i)))) \\ \text{s.t.} \quad & \text{ACC}(\text{LLM}(q_i; \mathcal{F}(\mathcal{X}'_K(q_i)))) \geq \theta. \end{aligned} \quad (1)$$

where  $\theta$  is a user-specified minimum acceptable accuracy. Overly conservative behavior (e.g., blanket refusals) may trivially reduce ASR but harms usability; our objective explicitly discourages such solutions by enforcing the ACC constraint.

## 4 Our Provable PRA-RAG

In RAG, we cast retrieval corruption as a *set perturbation* problem, where an adversary alters at most  $\varepsilon$  items in the Top- $K$  set, and propose a geometric, model-agnostic retrieval aggregation that operates on *combination embeddings*. As illustrated in Fig. 1, given a query  $q$ , we retrieve the Top- $K$  most relevant documents from a knowledge database, denoted as the set  $\mathcal{X} = \{p_1, \dots, p_K\}$ . In our approach, we set the number of retrieved

---

### Algorithm 1 Inference Procedure for Selecting Robust Retrieval Text

---

- 1: **Input:** A set of Top- $K$  retrieved texts  $\mathcal{X} = \{p_1, p_2, \dots, p_K\}$ ; subset size  $n$  (with constraint  $2n < K$ ).
  - 2: **Output:** A robustly aggregated retrieval text  $x^*$  for LLMs.
  - 3:  $\mathcal{S}_{\mathcal{X},n} = \{s_1, s_2, \dots, s_L\} \leftarrow \text{Comb}(\mathcal{X}, n)$ ,  $L = \binom{K}{n}$
  - 4: **for**  $i = 1$  to  $L$  **do**
  - 5:   embeddings  $\leftarrow [ \text{Embed}(p) \mid p \in s_i ]$
  - 6:    $v_i \leftarrow \text{Concat}(\text{embeddings})$
  - 7:    $\mathcal{V}_{\mathcal{X},n} \leftarrow \mathcal{V}_{\mathcal{X},n} \cup \{v_i\}$
  - 8: **end for**
  - 9: **for**  $i = 1$  to  $L$  **do**
  - 10:   Initialize similarity vector  $\text{sim}_i \in R^{L-1}$
  - 11:   **for**  $j = 1$  to  $L$  and  $j \neq i$  **do**
  - 12:      $d_{ij} \leftarrow \arccos(\cos(v_i, v_j))$
  - 13:     Store  $d_{ij}$  in  $\text{sim}_i$
  - 14:   **end for**
  - 15:   Sort  $\text{sim}_i$  in ascending order
  - 16:   Let  $m_i \leftarrow$  the  $k$ -th smallest value in  $\text{sim}_i$ ,  $k = \lceil \frac{L-1}{2} \rceil$
  - 17: **end for**
  - 18: Identify  $i^* = \arg \min_i m_i$ ,  $R = m_{i^*}$
  - 19: Compute  $x_{i^*} = \text{Aggregate}([ \text{Embed}(p) \mid p \in s_{i^*} ])$
  - 20: **return**  $x_{i^*}$  as input to the LLMs
- 

texts  $K$  slightly higher than the actual requirement to introduce diversity, thereby mitigating the impact of poisoned content. Next, we sample multiple subsets from the retrieved texts. Let  $\mathcal{S}_{\mathcal{X},n} = \{s_1, s_2, \dots, s_L\}$  denote the set of all possible subsets of size  $n$  drawn from the  $K$  retrieved documents, where  $L = \binom{K}{n}$ . Each document  $p_i$  in the subset is encoded as an embedding vector  $e_i \leftarrow \text{Embed}(p_i)$ ,  $e_i \in \mathbb{R}^d$ , yielding a set of embeddings in the  $d$ -dimensional vector space. We construct the corresponding vector set  $\mathcal{V}_{\mathcal{X},n} = \{v_1, v_2, \dots, v_L\}$  from the set  $\mathcal{S}$ , where each vector  $v_i$  is formed by concatenating the embedding vectors  $e_i$  of all texts in  $s_i$ .

We define the distance between two vectors  $u, v \in \mathcal{V}_{\mathcal{X},n}$  as the angular distance:  $d(u, v) = \arccos(\cos(u, v)) \in [0, \pi]$ , where  $\cos(u, v) = \frac{u \cdot v}{\|u\| \|v\|}$  denotes the cosine similarity. For a point  $v \in \mathcal{V}_{\mathcal{X},n}$  and radius  $R > 0$ , define the ball:

$$B(v, R) = \{u \in \mathcal{V}_{\mathcal{X},n} : d(u, v) \leq R\}. \quad (2)$$

Our method  $\mathcal{F}$  aims to identify the center of a minimum-radius ball that contains more than half of the points, representing the desired output of the robust aggregation process. Formally, the objective is defined as:

$$\mathcal{F}(\mathcal{X}) = \arg \min_z R \quad \text{s.t.} \quad \sum_{w \in \mathcal{V}_{\mathcal{X},n}} \mathbf{1}_{\mathcal{B}(z,R)}(w) \geq \left\lfloor \frac{\binom{K}{n}}{2} \right\rfloor + 1, \quad (3)$$

where  $\mathcal{X}$  denotes the input set,  $\mathcal{V}_{\mathcal{X},n}$  is the vector collection of all  $n$ -element subsets of  $\mathcal{X}$ ,  $\mathcal{B}(z, R)$  denotes a ball of radius  $R$  centered at  $z$ . The indicator function  $\mathbf{1}_{\mathcal{B}(z,R)}(w)$  is defined to be 1 if  $w \in \mathcal{B}(z, R)$ , and 0 otherwise. This formulation ensures that the selected center  $z$  lies within a ball that encompasses the majority of the subsets. Since  $S_{\mathcal{X},n}$  includes all possible combinations of the retrieved texts, the centroid  $z$  corresponding to the selected subset  $s$  effectively captures the meaning of the majority of texts in  $\mathcal{X}$ , thereby filtering out a small number of poisoned texts.

Finally, we compute a weighted average of the embedding vectors of the texts in the selected subset  $s_{i^*}$ , where the weights are determined by the similarity (e.g., cosine similarity) between each text  $p \in s_{i^*}$  and the query  $q$ :

$$x_{i^*} = \frac{\sum_{p \in s_{i^*}} \text{sim}(p, q) \cdot \text{Embed}(p)}{\sum_{p \in s_{i^*}} \text{sim}(p, q)}, \quad (4)$$

which helps prevent poisoned texts from evading filtration due to semantic similarity with clean texts, yielding the final robust aggregation result. More details are provided in Algorithm 1.

## 5 Computing the Certified Maximum Deviation

In the context of retrieval-based text poisoning attacks on RAG systems, we assume that the attacker successfully injects  $\varepsilon$  malicious documents into the final retrieved set  $\mathcal{X}$ , resulting in a perturbed set  $\mathcal{X}' = [\tilde{p}_1, \tilde{p}_2, \dots, \tilde{p}_\varepsilon] \cup [p_1, p_2, \dots, p_{K-\varepsilon}]$ , where each  $\tilde{p}_i$  represents a poisoned document inserted by the attacker, and each  $p_i$  denotes a clean document. In the absence of poisoning attacks, the center of the minimum-radius ball computed by Algorithm 1 is used as the robust output. Under poisoning attacks, the radius  $\hat{R}$  is enlarged to counter the shift caused by the  $\varepsilon$  poisoned texts, ensuring

that over half of the embeddings remain within the ball. Accordingly, the radius is approximately set to  $\hat{R} \leftarrow \mathcal{D}[k]$ , where  $k = \left\lfloor \frac{\binom{K}{n}}{2} \right\rfloor + \left( \binom{K}{n} - \binom{K-\varepsilon}{n} \right)$ , as shown in Algorithm 2. To ensure that the empirical radius  $\hat{R}$  satisfies the theoretical majority condition, Lemma 1 (Section A.4) guarantees that the majority ball contains a strict majority of clean combinations. Combined with the geometric separability assumption (Section A.5), this ensures the center is anchored to the clean subset. We denote the embedding shift of the final aggregated text caused by poisoned samples as  $d$ , where  $d = (1 + \beta)\hat{R}$ . The proposed metric  $d$  effectively quantifies the impact of poisoned texts on the retrieval results. A smaller  $d$  indicates lower sensitivity to perturbations, implying stronger robustness and greater resistance to successful attacks.

**Certifying the Upper Bound of Deviation  $d$ .** Let  $\mathcal{X}'$  be any corrupted version of  $\mathcal{X}$  obtained by changing at most  $\varepsilon$  passages, that is  $\|\mathcal{X}' - \mathcal{X}\|_0 \leq \varepsilon$ . Define the robustness radius  $R$  as the smallest value satisfying

$$\forall \mathcal{X}' \text{ s.t. } \|\mathcal{X}' - \mathcal{X}\|_0 \leq \varepsilon, \quad \sum_{w \in \mathcal{V}_{\mathcal{X}',n}} \mathbf{1}_{\mathcal{B}(\mathcal{F}(\mathcal{X}),R)}(w) \geq \left\lfloor \frac{\binom{K}{n}}{2} \right\rfloor + 1. \quad (5)$$

We present the following theorem:

**Theorem 1.** For  $\mathcal{X}'$  satisfying  $\|\mathcal{X}' - \mathcal{X}\|_0 \leq \varepsilon$ ,

$$d(\mathcal{F}(\mathcal{X}), \mathcal{F}(\mathcal{X}')) \leq 2R, \quad (6)$$

where the proof is provided in Section A.1. However, computing  $\mathcal{F}(\mathcal{X})$  exactly is computationally expensive in practice. Therefore, we use an approximation method from Center Smoothing (Kumar and Goldstein, 2021) to compute a  $\beta$ -MEB (Minimum Enclosing Ball) that contains the majority of elements in  $\mathcal{V}_{\mathcal{X},n}$ , denoted as:

$$\hat{\mathcal{F}}(\mathcal{X}) = \arg \min_{z \in \mathcal{V}_{\mathcal{X},n}} R \quad \text{s.t.} \quad \sum_{w \in \mathcal{V}_{\mathcal{X},n}} \mathbf{1}_{\mathcal{B}(z,R)}(w) \geq \left\lfloor \frac{\binom{K}{n}}{2} \right\rfloor + 1. \quad (7)$$

The radius of this approximation differs from the optimal radius by a multiplicative factor  $\beta$ . Similarly, the approximation of  $R$ , denoted as  $\hat{R}$ , is defined to satisfy:

---

**Algorithm 2** Certify Robustness under Poisoning

---

- 1: **Input:** A set of Top- $K$  retrieved texts  $\mathcal{X}$ ; subset size  $n$  (with  $2n < K$ ); maximum number of poisoned texts  $\varepsilon$  (with  $\binom{K}{n} < 2^{\binom{K-\varepsilon}{n}}$ ).
  - 2: **Output:** Certified maximum deviation  $d$ .
  - 3: Let  $\mathcal{S}_{\mathcal{X}',n} = \text{Comb}(\mathcal{X}', n)$
  - 4: Initialize  $\mathcal{V}_{\mathcal{X}',n} \leftarrow \emptyset$
  - 5: **for** each subset  $s_i \in \mathcal{S}_{\mathcal{X}',n}$  **do**
  - 6:    $v_i \leftarrow \text{Concat}([\text{Embed}(p) \mid p \in s_i])$
  - 7:    $\mathcal{V}_{\mathcal{X}',n} \leftarrow \mathcal{V}_{\mathcal{X}',n} \cup \{v_i\}$
  - 8: **end for**
  - 9: Let  $v_{\text{orig}}$  be the embedding vector of the subset selected by the Algorithm 1.
  - 10: Initialize an empty list of distances  $\mathcal{D} \leftarrow \emptyset$
  - 11: **for** each embedding  $v_i \in \mathcal{V}_{\mathcal{X}',n}$  **do**
  - 12:    $d_i \leftarrow \arccos(\cos(v_{\text{orig}}, v_i))$
  - 13:   Append  $d_i$  to  $\mathcal{D}$
  - 14: **end for**
  - 15: Sort  $\mathcal{D}$  in ascending order
  - 16: Let  $\hat{R} \leftarrow \mathcal{D}[k], k = \left\lfloor \frac{\binom{K}{n}}{2} \right\rfloor + \left( \binom{K}{n} - \binom{K-\varepsilon}{n} \right)$
  - 17: Let  $d = (1 + \beta) \hat{R}$
  - 18: **return**  $d$  as the certified maximum deviation
- 

$$\forall \mathcal{X}' \text{ s.t. } \|\mathcal{X}' - \mathcal{X}\|_0 \leq \varepsilon, \quad \sum_{w \in \mathcal{V}_{\mathcal{X}',n}} \mathbf{1}_{\mathcal{B}(\mathcal{F}(\mathcal{X}), \hat{R})}(w) \geq \left\lfloor \frac{\binom{K}{n}}{2} \right\rfloor + 1. \quad (8)$$

Then, we have the following theorem:

**Theorem 2.** For any  $\mathcal{X}'$  satisfying  $\|\mathcal{X}' - \mathcal{X}\|_0 \leq \varepsilon$ , we have:

$$d(\mathcal{F}(\mathcal{X}), \mathcal{F}(\mathcal{X}')) \leq (1 + \beta) \hat{R}, \quad (9)$$

where the proof is provided in Section A.2. We follow the practice in the work (Kumar and Goldstein, 2021) to use  $\beta = 2$  as an approximation when computing the minimum enclosing ball:

$$d(\mathcal{F}(\mathcal{X}), \mathcal{F}(\mathcal{X}')) \leq 3\hat{R}. \quad (10)$$

Note that while the theoretical formulation is based on the full set of  $L = \binom{K}{n}$  combinations, our framework naturally generalizes to a Monte Carlo sampling approach for larger values of  $K$  and  $n$ . In such cases, we sample  $m \ll L$  subsets, which provides a statistically bounded approximation of the optimal robustness (see Section A.3) and our experimental results further corroborate its effectiveness (see Section C.8).

## 6 Evaluation

### 6.1 Experimental Setup

In this section, we provide a detailed description of the experimental setup, with additional information and default configurations presented in Section B.1.

**Datasets.** In order to comprehensively evaluate the effectiveness of the proposed method, we utilize three question-answering datasets: Natural Questions (NQ) (Kwiatkowski et al., 2019), MS-MARCO (Bajaj et al., 2016), and HotpotQA (Yang et al., 2018).

**LLM.** We evaluate PRA-RAG with different LLMs, including Mistral-7B (Jiang et al., 2023), Llama3-8B, Vicuna-7B (Chiang et al., 2023), and GPT-3.5-Turbo (Brown et al., 2020).

**Attackers.** We conduct a systematic evaluation of PRA-RAG’s defense effectiveness against different types of poisoning attacks targeting the retrieved texts, following prior works (Xiang et al., 2024; Zhou et al., 2025): (1) Corpus Poisoning Attack: PoisonedRAG (Zou et al., 2024), (2) Adversarial Decoding: AdvDec (Zhang et al., 2025c), (3) Denial-Of-Service Attack: DoS Attack (Shafran et al., 2024), and (4) CorruptRAG (Zhang et al., 2025a). Further details of these attack methods are provided in Section B.3.

**Defenders.** We demonstrate the advantages of our approach by comparing it with representative state-of-the-art baselines, including RobustRAG (Xiang et al., 2024), InstructRAG (Wei et al., 2024), As-tuteRAG (Wang et al., 2025), TrustRAG (Zhou et al., 2025), and RAGForensics (Zhang et al., 2025b) with their default configurations.

**Evaluation Metrics.** We conduct a comprehensive evaluation of the proposed method using three metrics: (1) Adversarial Accuracy (**ACC**) assesses the system’s ability to produce correct answers when subjected to poisoning attacks; (2) Attack Success Rate (**ASR**) represents the proportion of incorrect answers generated due to poisoning attacks. We use GPT-4o to determine whether the LLM’s responses are correct or influenced by poisoning, using prompts are provided in the Section B.2; (3) Provable Average Deviation (**PAD**) provides a certified upper bound on semantic shifts in the aggregated retrieval representation. Lower PAD values indicate successful mitigation and preserved semantics, while higher values signal greater semantic corruption, making PAD a key metric for attack strength and robustness.

Table 1: The table presents the performance of PRA-RAG across different attack strategies and datasets under various LLMs (Top- $K = 8$ , subset size  $n = 3$ , poisoning rate = 20%).

Models	Metric	NQ					MS-MARCO					HotpotQA				
		PoisonedRAG	DoS Attack	AdvDec	CorruptRAG	Clean	PoisonedRAG	DoS Attack	AdvDec	CorruptRAG	Clean	PoisonedRAG	DoS Attack	AdvDec	CorruptRAG	Clean
Mistral-7B	ACC	0.62	0.57	0.55	0.55	0.60	0.63	0.65	0.61	0.59	0.57	0.37	0.33	0.39	0.41	0.38
	ASR	0.01	0.02	0.02	0.04	-	0.01	0.03	0.02	0.00	-	0.03	0.01	0.02	0.03	-
	PAD	1.23	1.43	1.30	1.21	0.79	1.37	1.46	1.51	1.39	0.90	1.33	1.55	1.43	1.31	0.83
Llama3-8B	ACC	0.47	0.49	0.52	0.49	0.50	0.69	0.71	0.68	0.67	0.63	0.30	0.33	0.31	0.30	0.29
	ASR	0.01	0.02	0.01	0.03	-	0.03	0.06	0.05	0.00	-	0.04	0.05	0.04	0.04	-
	PAD	0.95	1.04	0.97	0.93	0.60	1.09	1.19	1.06	1.11	0.75	1.06	1.17	1.07	1.04	0.65
Vicuna-7B	ACC	0.49	0.52	0.51	0.51	0.53	0.68	0.72	0.70	0.71	0.71	0.49	0.50	0.48	0.47	0.51
	ASR	0.01	0.03	0.02	0.00	-	0.01	0.02	0.03	0.00	-	0.00	0.02	0.02	0.04	-
	PAD	1.74	2.22	1.84	1.73	1.16	1.91	1.95	1.99	1.90	0.75	1.87	2.34	1.98	1.85	0.21
GPT-3.5-turbo	ACC	0.48	0.52	0.52	0.45	0.47	0.75	0.75	0.74	0.70	0.78	0.40	0.40	0.38	0.35	0.39
	ASR	0.00	0.02	0.00	0.03	-	0.01	0.00	0.02	0.00	-	0.07	0.08	0.02	0.06	-
	PAD	0.95	1.04	0.97	0.93	0.61	1.09	1.19	1.13	1.11	0.75	1.06	1.17	1.07	1.04	0.65

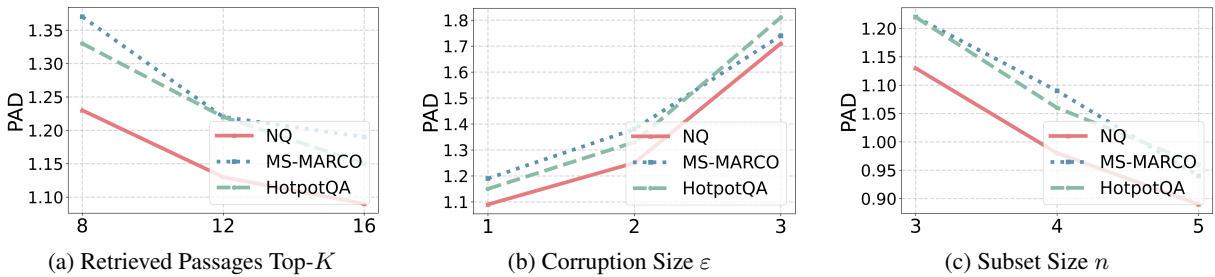


Figure 2: Impact of Retrieval Passages, Corruption Size and Subset Size on PAD scores for Mistral-7B.

## 6.2 Overall Results

**Main results of PRA-RAG.** Table 1 summarizes the performance of PRA-RAG under various poisoning attacks across multiple datasets and LLMs. The results demonstrate that PRA-RAG consistently achieves low ASR while maintaining high ACC, indicating strong robustness across all settings. For instance, on the MS-MARCO dataset with Mistral-7B, PRA-RAG attains an ACC of 62% under PoisonedRAG, with an ASR as low as 1%. Additionally, PAD effectively quantifies the semantic impact of adversarial retrievals: across models and tasks, lower PAD values correspond to lower ASR and higher ACC, supporting its utility as a reliable indicator of both robustness and attack severity. Additionally, in Section C.2, we compare the effects of various embedding extraction methods on the final model performance.

**PRA-RAG outperforms baselines.** Table 2 presents a comparative evaluation of PRA-RAG against several representative baselines under poisoning attacks and clean settings. Across all datasets and attack settings, PRA-RAG consistently achieves the best overall performance, attaining higher ACC while significantly lowering

the ASR. For instance, on the MSMARCO dataset with 20% poisoned retrievals, PRA-RAG maintains an ACC of 68% and suppresses ASR to just 1% with Vicuna-7B, outperforming RobustRAG and InstructRAG. Meanwhile, we observe that in the absence of poisoning attacks, the  $ACC_{clean}$  of our method remains largely unaffected. These results demonstrate that PRA-RAG not only offers stronger defense against poisoning attacks on retrieved texts but also preserves utility under normal conditions. Overall, existing methods either lack formal certification, trade usability for safety (over-refusal or over-filtering), or incur high cost with tight coupling to the LLM. PRA-RAG models retrieval poisoning as a set perturbation and performs robust aggregation on *combination embeddings* by selecting a majority-covering MEB center, yielding a certified PAD bound. Under the common regime, it consistently lowers ASR with low overhead, maintains competitive ACC, and withstands semantically-near poisoning.

## 6.3 Ablation Study

**Impact of retrieved passages Top- $K$ , corruption size  $\epsilon$  and subset size  $n$ .** The number of retrieved texts Top- $K$  and poisoned texts  $\epsilon$  significantly af-

Table 2: Performance of PRA-RAG and baselines under PoisonedRAG attacks (ACC, ASR) and clean settings (ACC<sub>clean</sub>). Results are based on the same experimental setup here for consistency and may differ from the original reports (Top- $K = 8$ , poisoning rate = 20%). Best results are in **bold**.

Models	Defense	NQ			MS-MARCO			HotpotQA		
		ACC <sub>clean</sub>	ACC	ASR	ACC <sub>clean</sub>	ACC	ASR	ACC <sub>clean</sub>	ACC	ASR
Mistral-7B	Vanilla RAG	0.61	0.26	0.46	0.58	0.33	0.18	0.55	0.17	0.48
	RobustRAG <sub>Keyword</sub>	<b>0.67</b>	0.66	0.09	0.61	<b>0.66</b>	0.12	<b>0.57</b>	<b>0.56</b>	0.24
	InstructRAG <sub>ICL</sub>	0.38	0.26	0.13	0.48	0.23	0.12	0.30	0.22	0.26
	ASTUTE RAG	0.48	0.47	0.05	0.58	0.48	0.05	0.35	0.35	0.09
	TrustRAG	0.65	<b>0.62</b>	0.03	<b>0.67</b>	0.64	0.08	0.50	0.45	0.10
	RAGForensics	0.41	0.41	0.07	0.62	0.65	0.03	0.45	0.46	0.10
	PRA-RAG <sub>cat</sub>	0.57	0.47	0.09	0.59	0.59	0.04	0.50	0.41	0.11
	PRA-RAG	0.60	<b>0.62</b>	<b>0.01</b>	0.57	0.63	<b>0.01</b>	0.49	0.37	<b>0.03</b>
	Vanilla RAG	<b>0.67</b>	0.31	0.45	0.73	0.35	0.33	<b>0.55</b>	0.23	0.43
Vicuna-7B	RobustRAG <sub>Keyword</sub>	0.48	0.40	0.08	0.59	0.40	0.07	0.44	0.31	0.10
	InstructRAG <sub>ICL</sub>	0.33	0.10	0.11	0.54	0.33	0.20	0.31	0.09	0.11
	ASTUTE RAG	0.50	0.44	0.19	<b>0.74</b>	0.63	0.14	0.34	0.26	0.30
	TrustRAG	0.51	0.49	0.06	0.62	0.55	0.06	0.36	0.30	0.07
	RAGForensics	0.60	<b>0.58</b>	0.03	0.68	<b>0.71</b>	0.04	0.51	<b>0.51</b>	0.06
	PRA-RAG <sub>cat</sub>	0.54	0.53	0.04	0.65	<b>0.71</b>	0.02	0.49	0.37	0.13
	PRA-RAG	0.53	0.49	<b>0.01</b>	0.63	0.68	<b>0.01</b>	0.50	0.49	<b>0.00</b>

fect both attack effectiveness and defense performance. Figure 2a illustrates the effect of varying the number of retrieved texts (Top- $K = 8, 12, 16$ ) on the PAD score of PRA-RAG, with  $\varepsilon = 1$  fixed. As Top- $K$  increases, the PAD score gradually declines, indicating enhanced robustness due to increased retrieval diversity. Figure 2b presents the opposite setting: with the number of retrieved texts fixed at Top- $K = 12$ , the PAD score increases as the number of poisoned texts rises from  $\varepsilon = 1$  to  $\varepsilon = 3$ . This trend indicates that injecting more poisoned texts into the retrieval results strengthens the attack, thereby diminishing the robustness of PRA-RAG. Figure 2c illustrates the impact of different subset sizes on the performance of PRA-RAG under the setting where the number of retrieved texts is Top- $K = 12$  and the number of poisoned texts is  $\varepsilon = 1$ . We evaluate cases where each combination contains 3, 4, or 5 texts. The results show a clear downward trend in PAD scores as the number of texts per subset increases. This indicates that incorporating more clean texts within each combination helps dilute the influence of poisoned content and improves robustness. Further experimental data and analyses can be found in Section C.3.

**Comparison of aggregation strategies.** We compare two aggregation methods for the selected subset: a baseline that concatenates texts as input to the LLM (PRA-RAG<sub>cat</sub>), and our method (PRA-RAG), which computes a weighted average of embeddings.

Table 3: The table presents the average response latency of different RAG methods across datasets.

Dataset	ASTUTE RAG	InstructRAG <sub>ICL</sub>	RobustRAG <sub>Keyword</sub>	PRA-RAG
NQ	13.68s	6.89s	27.01s	<b>5.98s</b>
MS-MARCO	13.84s	7.15s	26.38s	<b>6.20s</b>
HotpotQA	16.78s	7.01s	25.87s	<b>6.17s</b>

As shown in Table 2, our approach achieves higher ACC and significantly reduces ASR by better mitigating the impact of poisoned texts. In Section C.4, we further compare and analyze the impact of different distance metrics used to compute  $d$  on the performance of our method.

## 6.4 Efficiency

In real-world applications, response time is critical to the user experience of RAG systems. However, most existing defense strategies often incur varying levels of inference overhead. In our experiments, we evaluated the average response latency over 100 query examples to compare different defense methods. As shown in Table 3, our proposed PRA-RAG achieves the lowest inference latency, demonstrating its practical efficiency while preserving strong robustness. In our framework, both the number of retrieved texts (Top- $K$ ) and the subset size have a significant impact on system efficiency. Therefore, we conduct a more in-depth analysis and empirical evaluation of the computational overhead in Section C.5. Additionally, Section C.8 presents the computational overhead and analysis under Monte Carlo sampling.

## 7 Conclusion

This paper presents PRA-RAG, a provably robust aggregation method for defending against retrieval-level poisoning attacks in RAG systems. We theoretically derive an upper bound on the semantic deviation introduced by poisoned retrievals and propose Provable Average Deviation (PAD) as a unified metric for evaluating both robustness and attack strength. Unlike traditional defenses that depend on model outputs, PRA-RAG ensures semantic stability directly in the embedding space, providing a principled, model-agnostic robustness guarantee. Extensive experiments show that PRA-RAG significantly enhances robustness against poisoning attacks across diverse datasets and language models, consistently lowering attack success rates while preserving high answer accuracy.

## 590 Limitations.

591 Our work has the following limitations:

- 592 • This work selects the most robust subset of the  
593 retrieved texts as the final input to the LLM  
594 for answer generation. Although the selected  
595 subset may not always produce the optimal  
596 response from the LLM, as shown in our ex-  
597 periments, the impact on model utility remains  
598 limited while significantly enhancing robust-  
599 ness against poisoning attacks.
- 600 • Our approach enhances robustness by con-  
601 structing subsets, which increases computa-  
602 tional and response-time overhead. However,  
603 experiments show that as the number and size  
604 of subsets grow, the defense against poison-  
605 ing attacks significantly improves. In practice,  
606 users can flexibly balance efficiency and secu-  
607 rity based on their needs.
- 608 • The effectiveness of our approach relies on the  
609 number of poisoned texts remaining below a  
610 certain threshold. Although an attacker could  
611 potentially bypass the defense through large-  
612 scale poisoning, this incurs substantial costs  
613 and significantly reduces the quality of the  
614 context, making it difficult for even humans  
615 to provide correct answers.

## 616 Ethics Statement

617 The goal of this work is to defend against  
618 retrieval-based poisoning attacks in RAG systems.  
619 All data used in this study is publicly available, en-  
620 suring no additional privacy concerns. The source  
621 code and software will be released as open-source.  
622 While this openness may expose the system to adap-  
623 tive attacks, our approach can be further strength-  
624 ened by incorporating additional internal and exter-  
625 nal information. Overall, we believe our method  
626 contributes to advancing the secure deployment of  
627 RAG systems.

## 628 References

629 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
630 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
631 Diogo Almeida, Janko Altenschmidt, Sam Altman,  
632 Shyamal Anadkat, and 1 others. 2023. Gpt-4 techni-  
633 cal report. *arXiv preprint arXiv:2303.08774*.

634 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and  
635 Hannaneh Hajishirzi. 2024. Self-rag: Learning to re-  
636 trieve, generate, and critique through self-reflection.

Mihai Badoiu and Kenneth L Clarkson. 2003. Smaller  
637 core-sets for balls. In *SODA*, volume 3, pages 801–  
638 802. 639

Jinheon Baek, Soyeong Jeong, Minki Kang, Jong C  
640 Park, and Sung Hwang. 2023. Knowledge-  
641 augmented language model verification. In *Proceed-*  
642 *ings of the 2023 Conference on Empirical Methods*  
643 *in Natural Language Processing*, pages 1720–1736. 644

Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng,  
645 Jianfeng Gao, Xiaodong Liu, Rangan Majumder, An-  
646 drew McNamara, Bhaskar Mitra, Tri Nguyen, and  
647 1 others. 2016. Ms marco: A human generated ma-  
648 chine reading comprehension dataset. *arXiv preprint*  
649 *arXiv:1611.09268*. 650

BBC. 2024. Glue pizza and eat rocks: Google ai  
651 search errors go viral. [https://www.bbc.co.uk/](https://www.bbc.co.uk/news/articles/cd11gzejgz4o)  
652 [news/articles/cd11gzejgz4o](https://www.bbc.co.uk/news/articles/cd11gzejgz4o). 653

Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
654 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
655 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
656 Askell, and 1 others. 2020. Language models are  
657 few-shot learners. *Advances in neural information*  
658 *processing systems*, 33:1877–1901. 659

Zirui Cheng, Jikai Sun, Anjun Gao, Yueyang Quan,  
660 Zhuqing Liu, Xiaohua Hu, and Minghong Fang. 2025.  
661 Secure retrieval-augmented generation against poi-  
662 soning attacks. *arXiv preprint arXiv:2510.25025*. 663

Wei-Lin Chiang, Zhuohan Li, Ziqing Lin, Ying Sheng,  
664 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan  
665 Zhuang, Yonghao Zhuang, Joseph E Gonzalez, and  
666 1 others. 2023. Vicuna: An open-source chatbot  
667 impressing gpt-4 with 90%\* chatgpt quality. *See*  
668 <https://vicuna.lmsys.org> (accessed 14 April 2023),  
669 2(3):6. 670

Sukmin Cho, Jeongyeon Seo, Soyeong Jeong, and  
671 Jong C Park. 2023. Improving zero-shot reader by  
672 reducing distractions from irrelevant documents in  
673 open-domain question answering. In *Findings of the*  
674 *Association for Computational Linguistics: EMNLP*  
675 *2023*, pages 3145–3157. 676

Google. 2024. Generative ai in search:  
677 Let google do the searching for you.  
678 [https://blog.google/products/search/](https://blog.google/products/search/generative-ai-google-search-may-2024/)  
679 [generative-ai-google-search-may-2024/](https://blog.google/products/search/generative-ai-google-search-may-2024/). 680

Kai Greshake, Sahar Abdelnabi, Shailesh Mishra,  
681 Christoph Endres, Thorsten Holz, and Mario Fritz.  
682 2023. Not what you’ve signed up for: Compromis-  
683 ing real-world llm-integrated applications with indi-  
684 rect prompt injection. In *Proceedings of the 16th*  
685 *ACM Workshop on Artificial Intelligence and Secu-*  
686 *rity*, pages 79–90. 687

Wassily Hoeffding. 1963. Probability inequalities for  
688 sums of bounded random variables. *Journal of the*  
689 *American statistical association*, 58(301):13–30. 690



Baolei Zhang, Haoran Xin, Minghong Fang, Zhuqing Liu, Biao Yi, Tong Li, and Zheli Liu. 2025b. Traceback of poisoning attacks to retrieval-augmented generation. In *Proceedings of the ACM on Web Conference 2025*, pages 2085–2097.

Collin Zhang, Tingwei Zhang, and Vitaly Shmatikov. 2025c. Adversarial decoding: Generating readable documents for adversarial objectives. *Preprint*, arXiv:2410.02163.

Huichi Zhou, Kin-Hei Lee, Zhonghao Zhan, Yue Chen, and Zhenhao Li. 2025. Trustrag: Enhancing robustness and trustworthiness in rag. *arXiv preprint arXiv:2501.00879*.

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

## A Theoretical Analysis and Proofs

### A.1 Proof of Theorem 1

Consider two balls  $\mathcal{B}(\mathcal{F}(\mathcal{X}'), r^*(\mathcal{X}'))$  and  $\mathcal{B}(\mathcal{F}(\mathcal{X}), R)$ . According to the definitions of  $R$  and  $r^*(\mathcal{X}')$ , these two balls must share at least one common element, denoted as  $w^*$ . Since the defined distance based on vector angles satisfies the triangle inequality, we have:

$$\begin{aligned} d(\mathcal{F}(\mathcal{X}), \mathcal{F}(\mathcal{X}')) &\leq d(\mathcal{F}(\mathcal{X}), w^*) + d(\mathcal{F}(\mathcal{X}'), w^*) \\ &\leq R + r^*(\mathcal{X}'). \end{aligned} \quad (11)$$

Since the ball  $\mathcal{B}(\mathcal{F}(\mathcal{X}), R)$  contains at least a majority of elements in  $\mathcal{V}_{\mathcal{X}',n}$ , the smallest radius  $r^*(\mathcal{X}')$  of a ball that contains a majority of elements in  $\mathcal{V}_{\mathcal{X}',n}$  must not be greater than  $R$ , i.e.,  $r^*(\mathcal{X}') \leq R$ . Therefore,

$$d(\mathcal{F}(\mathcal{X}), \mathcal{F}(\mathcal{X}')) \leq 2R. \quad (12)$$

### A.2 Proof of Theorem 2

Since the radius of the  $\beta$ -MEB differs from the optimal radius by a factor  $\beta$ , i.e., the radius of the  $\beta$ -MEB over  $\mathcal{V}_{\mathcal{X},n}$  is  $\beta r^*(\mathcal{X}')$ , by the same reasoning as in Theorem 1, we have:

$$\begin{aligned} d(\widehat{\mathcal{F}}(\mathcal{X}), \widehat{\mathcal{F}}(\mathcal{X}')) &\leq d(\widehat{\mathcal{F}}(\mathcal{X}), w^*) + d(\widehat{\mathcal{F}}(\mathcal{X}'), w^*) \\ &\leq \widehat{R} + \beta r^*(\mathcal{X}') \leq (1 + \beta)\widehat{R}. \end{aligned} \quad (13)$$

### A.3 Scalability and Approximation Guarantee via Monte Carlo Sampling

To address the rapid increase in computational cost as  $K$  and  $n$  grow, we introduce a Monte Carlo sampling strategy. Instead of enumerating all  $L = \binom{K}{n}$

combinations, we randomly sample  $m$  subsets to estimate the robust aggregation. Here, we prove that the radius  $R^*$  derived from this sampling strategy provides a theoretically bounded approximation of the optimal radius  $R$  derived from the full set.

**Theorem 3 (Approximation Guarantee of Sampling).** Let  $\mathcal{V}_{\mathcal{X},n}$  be the set of all embedding vectors generated from  $\binom{K}{n}$  combinations, and let  $R$  be the radius of the Minimum Enclosing Ball (MEB) covering the majority ( $> 50\%$ ) of  $\mathcal{V}_{\mathcal{X},n}$ . Let  $\mathcal{S} \subset \mathcal{V}_{\mathcal{X},n}$  be a random sample of size  $m$ , and let  $R^*$  be the radius of the MEB covering the majority of  $\mathcal{S}$ . For any  $\tau > 0$ , if the sample size satisfies  $m \gg \lceil 2/\tau \rceil$ , then with high probability:

$$R^* \leq R \leq (1 + \tau)R^*. \quad (14)$$

*Proof.* The proof relies on establishing the statistical consistency of the majority, followed by the geometric bounds.

**Majority Preservation (Statistical Consistency):** First, we must ensure that the majority semantics are preserved during sampling. Let  $p$  be the proportion of clean subsets in the full set  $\mathcal{V}_{\mathcal{X},n}$ , and  $\hat{p}$  be the proportion in the sample  $\mathcal{S}$ . According to Hoeffding’s Inequality (Hoeffding, 1963), the probability that the sampled proportion deviates from the true proportion is bounded by:

$$\mathbb{P}(|\hat{p} - p| \geq \delta) \leq 2e^{-2m\delta^2}. \quad (15)$$

For a sample size of  $m = 200$  and error margin  $\delta = 0.1$ , the probability of failure is negligible ( $< 0.04$ ). This indicates that under a mild assumption that the clean subsets constitute a sufficient majority in the full set (i.e.,  $p \geq 0.5 + \delta$ ), they are guaranteed to remain the majority ( $\hat{p} > 0.5$ ) in the sampled subset  $\mathcal{S}$  with high probability. As a result, the optimization target (i.e., the clean cluster) remains consistent between the full set and the sample. Furthermore, since our sampling is uniformly performed over the entire set, the sampled subset  $\mathcal{S}$  is guaranteed with high probability to follow the same distribution as the full set  $\mathcal{V}_{\mathcal{X},n}$ .

**Lower Bound ( $R^* \leq R$ ):** Let  $B$  be the optimal ball with radius  $R$  that covers the majority of the full set  $\mathcal{V}_{\mathcal{X},n}$ . Due to the majority preservation property established above,  $B$  also covers the majority of the points in the random sample  $\mathcal{S}$  with high probability. Since  $R$  is defined as the *minimum* radius necessary to cover the majority of  $\mathcal{S}$  (the optimal solution for the sample), and  $R$  is a

valid radius that satisfies this condition (a feasible solution for the sample), it follows from the definition of optimality that:

$$R^* \leq R. \quad (16)$$

**Upper Bound** ( $R \leq (1 + \tau)R^*$ ): We invoke the *Core-Set Theory* for Minimum Enclosing Balls (Badoiu and Clarkson, 2003). The theorem establishes that for any geometric set, there exists a core-set of size  $\lceil 2/\tau \rceil$  such that the MEB of the core-set approximates the true MEB within a factor of  $(1 + \tau)$ . In our context, we apply this theorem specifically to the clean geometric cluster (which constitutes the majority of  $\mathcal{V}_{\mathcal{X},n}$ ). Since our sample size  $m$  is sufficiently large ( $m \gg \lceil 2/\tau \rceil$ ), the random sample  $\mathcal{S}$  captures the core-set of this clean cluster with high probability. Because the sample  $\mathcal{S}$  contains the core-set, the radius  $R^*$  computed from  $\mathcal{S}$  is at least as large as the radius of the core-set. Consequently, expanding  $R^*$  by  $(1 + \tau)$  is sufficient to cover the entire clean cluster, which corresponds to the majority of the full population:

$$R \leq (1 + \tau)R^*. \quad (17)$$

**Conclusion.** Combining the probabilistic consistency and geometric bounds, we conclude that Monte Carlo sampling with  $m \approx 200$  and  $\tau = 0.5$  yields a  $(1 + \tau)$ -approximation of the true robustness metric PAD, effectively decoupling the computational complexity from the combinatorial growth of  $K$ , reducing it to a constant factor determined by the sample size  $m$ .

#### A.4 Guarantee of Clean Majority in the Certified Ball

**Lemma 1.** Let  $\mathcal{V}$  be the set of all  $N = \binom{K}{n}$  subset embeddings derived from the retrieved documents. Let  $\mathcal{C} \subset \mathcal{V}$  denote the set of clean embeddings and  $\mathcal{P} \subset \mathcal{V}$  denote the set of poisoned embeddings, such that  $|\mathcal{P}| \leq N_{adv} = \binom{K}{n} - \binom{K-\varepsilon}{n}$ . If the certified radius  $\hat{R}$  is determined by the distance value at the  $k$ -th index (0-based) of the sorted distance vector to the geometric center, with  $k = \lfloor N/2 \rfloor + N_{adv}$ , then the ball  $\mathcal{B}(z, \hat{R})$  is guaranteed to strictly contain a majority of clean embeddings, i.e.,  $|\mathcal{B}(z, \hat{R}) \cap \mathcal{C}| \geq \lfloor N/2 \rfloor + 1$ .

**Proof.** Let  $\mathcal{S}$  denote the set of embeddings enclosed by the ball  $\mathcal{B}(z, \hat{R})$ . Since the radius  $\hat{R}$  corresponds to the distance at the  $k$ -th index (where index 0 represents the center itself), the set  $\mathcal{S}$  includes the first  $k + 1$  smallest distance embeddings.

Thus, by definition,  $|\mathcal{S}| = k + 1$ . The set  $\mathcal{S}$  is composed of disjoint clean and poisoned subsets:  $\mathcal{S} = (\mathcal{S} \cap \mathcal{C}) \cup (\mathcal{S} \cap \mathcal{P})$ . We seek a lower bound on the number of clean samples  $|\mathcal{S} \cap \mathcal{C}|$ . According to set theory:

$$|\mathcal{S} \cap \mathcal{C}| = |\mathcal{S}| - |\mathcal{S} \cap \mathcal{P}| \quad (18)$$

In the worst-case adversarial scenario, the adversary optimizes the poisoned embeddings to be geometrically closest to the center to occupy the top ranks. However, the total number of poisoned embeddings is strictly bounded by  $N_{adv}$ . Thus, for any selected subset,  $|\mathcal{S} \cap \mathcal{P}| \leq N_{adv}$  always holds. Substituting the cardinality  $|\mathcal{S}| = k + 1$  and the value of  $k$ :

$$\begin{aligned} |\mathcal{S} \cap \mathcal{C}| &\geq (k + 1) - N_{adv} \\ &= (\lfloor N/2 \rfloor + N_{adv} + 1) - N_{adv} \\ &= \lfloor N/2 \rfloor + 1 \end{aligned} \quad (19)$$

Therefore, the ball  $\mathcal{B}(z, \hat{R})$  necessarily contains at least  $\lfloor N/2 \rfloor + 1$  clean embeddings. Since  $N = |\mathcal{V}|$ , this count represents a strict majority of the total universe of combinations, explicitly satisfying the condition required for Theorem 1.

#### A.5 Geometric Separability Assumption

We assume a mild geometric separability condition to connect majority coverage over combination embeddings to clean subset selection. Specifically, there exists a center  $z^*$  and a radius  $R_c$  such that a strict majority of fully clean combination embeddings lie within the ball  $\mathcal{B}(z^*, R_c)$ , while any combination containing at least one poisoned passage lies outside  $\mathcal{B}(z^*, R_c + \Delta)$  for some  $\Delta > 0$ .

The separability condition is necessary to ensure that a ball containing a strict majority of clean embeddings is geometrically anchored to the clean subset, thereby enabling the selection of clean contexts through geometric aggregation. When the margin  $\Delta = 0$ , clean and poisoned combination embeddings become geometrically indistinguishable. In such cases, no geometry-based majority certification method can guarantee the selection of a fully clean subset, as the observed information no longer provides distinguishable features.

## B Details

### B.1 Details of Experiment Setup

**Dataset.** Detailed statistics are shown in Table 4.

**Default Setting.** Unless otherwise specified, our

Table 4: Statics of datasets.

Datasets	Texts in knowledge database	Questions
NQ (Kwiatkowski et al., 2019)	2,681,468	3,452
MS-MARCO (Bajaj et al., 2016)	5,233,329	7,405
HotpotQA (Yang et al., 2018)	8,841,823	6,980

experiments use the following default settings: the Contriever (Izacard et al., 2021) retriever is employed to select the Top- $K = 8$  texts most similar to the query from the MS-MARCO dataset, with poisoned texts generated according to the PoisonedRAG (Zou et al., 2024) method constituting 20% (rounded down) of the retrieval results. Subsequently, multiple subsets are sampled from the retrieved texts, each containing  $n = 3$  texts. Open-source LLMs can directly obtain embeddings through their embedding layers and generate the final robust aggregated representation. For black-box models (e.g., GPT-3.5-turbo), we use Llama3-8B in our experiments to extract embeddings and decode the robust aggregation results into text to serve as context. Finally, the Mistral-7B language model generates responses based on these retrieved texts. Each result is averaged over 10 runs under a consistent setup.

**Implementation and Resources.** We report the mean performance over 10 independent runs and assess the statistical significance of improvements or degradations using the Wilcoxon signed-rank test at a significance level of  $p \leq 0.05$ . We conduct experiments on a server with 64 AMD EPYC 9654 CPUs at 2.40–3.70 GHz, 512 GB of DDR5 RAM (assumed based on high-core-count server standards), and four NVIDIA RTX A100 GPUs, each with 80 GB GDDR6 memory.

## B.2 Prompt

You are a helpful assistant. The user has provided a query along with relevant context information. Use this context to answer the question briefly and clearly. If you cannot find the answer to the question, respond with "I don't know."  
 Contexts: [context]  
 Query: [question]  
 Answer:

Figure 3: The prompt used in RAG to make an LLM generate an answer based on the retrieved texts.

You are given a question and two answers. Determine whether the two answers are semantically consistent with each other, whether they express the same stance, conclusion, or outcome regarding the question. If they are consistent, output "yes". If they are inconsistent or contradictory, output "no".  
 Query: [question]  
 Answers: [correct/incorrect answer & response of LLM]  
 Answer:

Figure 4: The prompt used in RAG to determine whether the LLM’s response is correct or corresponds to the attacker’s target response.

## B.3 Poisoned Texts Generation

To comprehensively evaluate the robustness of our approach against different poisoned text generation strategies, we adopt three representative strategies, including PoisonedRAG (Zou et al., 2024), Adversarial Decoding (Zhang et al., 2025c), Denial-Of-Service Attack (Shafran et al., 2024), and CorruptRAG (Zhang et al., 2025a). In PoisonedRAG (Zou et al., 2024) method, the attacker begins by selecting a target question and its corresponding incorrect answer, then crafts poisoned texts to satisfy two key requirements: (1) being retrievable by the retriever and (2) successfully misleading the language model into producing the incorrect answer. Adversarial Decoding (Zhang et al., 2025c) is a token-level beam search framework that integrates continuous, task-specific scorers into standard decoding, enabling the generation of fluent, low-detectability texts that jointly optimize retrieval similarity and adversarial generation objectives. Denial-Of-Service Attack (Shafran et al., 2024) inserts a single blocker document consisting of a retrieval component to ensure retrieval and a jamming component to trigger refusal responses, crafted through instruction injection, oracle generation, or black-box optimization. CorruptRAG (Zhang et al., 2025a) can efficiently mislead the model into generating targeted false content by injecting only a single poisoned text that simultaneously includes the target query, which ensures it is preferentially retrieved, and an adversarial prompting template, which exploits the LLM’s generation bias to label the correct answer as outdated and to fabricate a supposedly latest but incorrect answer. In our experiments, we adopt the default parameters of these methods to generate the corresponding poisoned texts, ensuring effective poisoning.

## C Additional Experimental Results

### C.1 Distribution of Different Combinations

Figure 5 illustrates the distributional differences between combinations containing poisoned texts and clean combinations. It is evident that the poisoned texts, intentionally crafted to induce attacker-desired responses, exhibit a significant shift in their embedding vectors compared to clean texts.

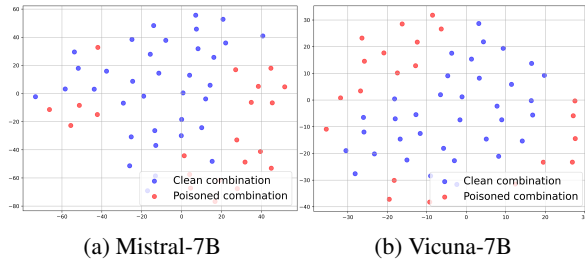


Figure 5: Embedding Distribution of Retrieved Text Combinations (Red: With Poisoned Texts; Blue: Without Poisoned Texts)

### C.2 Different Embedding Extraction Methods

Table 5 compares the impact of different embedding acquisition methods on the final performance of PRA-RAG. In our experiments with Mistral-7B, LLaMA-3-8B, and GPT-3.5-turbo, embeddings are obtained either from each model’s native embedding module or from the Contriever retriever. While the results show minimal variation in ACC and ASR, they reveal substantial discrepancies in PAD, reflecting the differences in how embeddings are computed. This highlights the importance of standardizing the embedding source in practice to ensure consistent robustness. This further reflects that our method is robust and generalizable across different embedding models.

Table 5: The table summarizes the performance of PRA-RAG using different embedding extraction methods under GPT-3.5-turbo.

Models	Methods	NQ			MS-MARCO			HotpotQA		
		ACC	ASR	PAD	ACC	ASR	PAD	ACC	ASR	PAD
Mistral-7B	Contriever	0.55	0.03	2.52	0.62	0.02	2.62	0.44	0.08	2.71
	Mistral-7B	0.62	0.01	1.23	0.63	0.01	1.37	0.37	0.03	1.33
Llama3-8B	Contriever	0.52	0.04	2.52	0.65	0.00	2.62	0.33	0.06	2.71
	Llama3-8B	0.47	0.01	0.95	0.69	0.03	1.09	0.30	0.04	1.06
GPT-3.5-turbo	Contriever	0.47	0.02	2.52	0.78	0.01	2.62	0.32	0.09	2.71
	Llama3-8B	0.48	0.00	0.95	0.75	0.00	1.09	0.40	0.07	1.06

### C.3 Impact of Different Parameters

**Impact of Retrieved Passages Top- $K$ .** Table 6 presents the impact of varying the number of retrieved texts on the performance of PRA-RAG under the setting where the number of poisoned texts is 1 (10%) and the subset size is 3. The experimental results show that the changes in ACC and ASR vary across different large language models and datasets, lacking a consistent trend. However, the PAD metric exhibits a consistently decreasing trend. This is because increasing the number of retrieved texts reduces the proportion of poisoned texts, thereby diminishing their overall influence.

Table 6: The table shows the performance of PRA-RAG with different numbers of retrieved texts.

Models	Top- $K$	NQ			MS-MARCO			HotpotQA		
		ACC	ASR	PAD	ACC	ASR	PAD	ACC	ASR	PAD
Mistral-7B	8	0.62	0.01	1.23	0.63	0.01	1.37	0.37	0.03	1.33
	12	0.54	0.00	1.13	0.66	0.00	1.22	0.40	0.02	1.22
	16	0.53	0.00	1.09	0.53	0.00	1.19	0.40	0.02	1.15
Vicuna-7B	8	0.49	0.01	1.74	0.68	0.01	1.91	0.49	0.00	1.87
	12	0.53	0.00	1.66	0.68	0.00	1.73	0.52	0.01	1.76
	16	0.52	0.00	1.59	0.70	0.00	1.69	0.51	0.01	1.68

**Impact of corruption size  $\epsilon$ .** Table 7 reports the effect of the number of poisoned texts on the robustness of PRA-RAG with the number of retrieved texts Top- $K = 12$ . The results show that as the number of poisoned texts increases, both ASR and PAD rise across models and datasets, indicating stronger attacks and reduced robustness. More experimental results and analysis under different poisoning ratios are provided in Section C.7

Table 7: The table shows the performance of PRA-RAG with different numbers of poisoned texts (Top- $K = 12$ , subset size = 3). Here,  $\epsilon = 1, 2, 3$  are obtained by multiplying the poisoning rates (10%, 20%, 30%) by 12 and taking the floor.

Models	Corruption Size	NQ			MS-MARCO			HotpotQA		
		ACC	ASR	PAD	ACC	ASR	PAD	ACC	ASR	PAD
Mistral-7B	1	0.53	0.00	1.09	0.53	0.00	1.19	0.40	0.02	1.15
	2	0.51	0.00	1.25	0.56	0.00	1.38	0.38	0.05	1.33
	3	0.49	0.04	1.71	0.53	0.03	1.74	0.33	0.09	1.81
Vicuna-7B	1	0.52	0.00	1.59	0.70	0.00	1.69	0.51	0.01	1.68
	2	0.48	0.00	1.79	0.71	0.00	1.93	0.50	0.02	1.89
	3	0.50	0.01	2.23	0.70	0.02	2.45	0.50	0.04	2.43

**Impact of subset size  $n$ .** Table 8 presents the impact of the number of texts per sampled subset

Table 8: The table shows the performance of PRA-RAG with different subset size.

Models	Subset Size	NQ			MS-MARCO			HotpotQA		
		ACC	ASR	PAD	ACC	ASR	PAD	ACC	ASR	PAD
Mistral-7B	3	0.54	0.00	1.13	0.66	0.00	1.22	0.40	0.02	1.22
	4	0.60	0.01	0.98	0.63	0.00	1.09	0.44	0.01	1.06
	5	0.63	0.00	0.89	0.71	0.01	0.94	0.43	0.01	0.96
Vicuna-7B	3	0.53	0.00	1.66	0.68	0.00	1.73	0.52	0.01	1.76
	4	0.55	0.01	1.43	0.67	0.00	1.53	0.54	0.00	1.53
	5	0.58	0.00	1.27	0.69	0.01	1.38	0.58	0.01	1.37

on the results when the total number of retrieved texts is fixed at  $Top-K = 12$ . The experimental results show that as the number of texts included in each subset increases, the context provided to the LLM becomes richer, leading to an upward trend in ACC but a downward trend in PAD.

#### C.4 Comparison of different distance metrics

In the theoretical framework of PRA-RAG, we adopt *angular distance* to ensure strict metric properties. In implementation, this aligns with maximizing *cosine similarity*, which focuses on vector direction. To evaluate the impact of this metric choice, we compare it against *Manhattan distance* and *Standard Euclidean distance*. As shown in Table 9, the angular distance-based method consistently outperforms the standard Euclidean baseline across all evaluation metrics. This confirms that capturing semantic alignment (via angular distance) is more effective than measuring absolute positional differences for this task.

Table 9: The table presents the performance of PRA-RAG with different distance metrics.

Distance Metrics	NQ		MS-MARCO		HotpotQA	
	ACC	ASR	ACC	ASR	ACC	ASR
<i>Manhattan distance</i>	0.58	0.01	0.59	0.02	0.36	0.12
<i>Euclidean distance</i>	0.61	0.03	0.60	0.02	0.36	0.11
<i>angular distance</i>	<b>0.62</b>	0.01	<b>0.63</b>	<b>0.01</b>	<b>0.37</b>	<b>0.03</b>

#### C.5 Analysis of Computational Overhead

Table 10 compares the computational overhead under different numbers of retrieved texts, with the subset size  $n = 3$ . The results show that as  $Top-K$  increases, the number of combinations grows rapidly, leading to a near-linear increase in computation time. Specifically, when  $Top-K$  is set to 8, 12, and 16, the average generation times are 5.98

Table 10: Computational overhead of PRA-RAG under different  $Top-K$  retrieval settings, with subset size fixed at 3 and a poisoning rate of 20%.

Models	Top- $K$	NQ	MS-MARCO	HotpotQA
Mistral-7B	8	5.98s	6.20s	6.17s
	12	18.47s	24.88s	22.14s
	16	68.05s	65.48s	72.16s
Vicuna-7B	8	3.37s	2.91s	3.02s
	12	12.13s	11.58s	12.23s
	16	64.93s	67.13s	59.25s

Table 11: Computational overhead of PRA-RAG with  $Top-K = 12$  retrieval under a poisoning rate of 10%, evaluated across different subset sizes.

Model	Subset Size	NQ	MS-MARCO	HotpotQA
Mistral-7B	3	18.71s	18.10s	17.76s
	4	55.25s	55.30s	57.89s
	5	125.55s	131.37s	124.96s
Vicuna-7B	3	11.88s	12.16s	11.84s
	4	47.92s	46.12s	52.40s
	5	116.09s	121.47s	117.85s

s, 18.47 s, and 68.05 s, respectively. Table 11 compares the computational overhead under a fixed  $Top-K$  retrieval size of 12 with varying subset sizes. The results show that the computation time no longer correlates linearly with the number of combinations, as changes in subset size also affect the overall time required for the final RAG response generation. Both the number of retrieved texts ( $Top-K$ ) and the subset size ( $n$ ) have a significant impact on the response time of our approach. Nevertheless, as shown in Table 3, our method still outperforms existing baselines in terms of efficiency under certain settings (e.g.,  $Top-K = 8$ ,  $n = 3$ ), which represent typical real-world scenarios. Although increasing  $Top-K$  and  $n$  leads to higher computational overhead, it also enhances system security. As shown in Figure 2, the PAD value decreases as  $Top-K$  and  $n$  increase, indicating higher security. Overall, users can flexibly balance efficiency and security by selecting appropriate values of  $Top-K$  and  $n$  based on practical needs.

#### C.6 Robustness of PRA-RAG Against Adaptive Attacks

We extended our evaluation to an adaptive attack scenario in which the attacker generates poisoned texts with high semantic similarity to the correct documents. We utilized PoisonedRAG (Zou et al., 2024) with additional semantic constraints to align the poisoned texts with the ground truth in both

Table 12: Performance of PRA-RAG under adaptive attack scenarios.

Models	Metric	NQ	MS-MARCO	HotpotQA
Mistral-7B	ACC	0.60	0.57	0.38
	ASR	0.00	0.00	0.04
	PAD	1.23	1.36	1.30
Vicuna-7B	ACC	0.48	0.68	0.52
	ASR	0.01	0.02	0.05
	PAD	1.73	1.88	1.83

content and format while maintaining attack effectiveness. With  $Top-K = 8$ ,  $n = 3$ , and a 20% poisoning rate, the results in Table 12 confirm the robustness of our method against adaptive attacks. For example, the Vicuna-7B model on the NQ dataset yielded an ACC of 0.48 and an ASR of 0.01, closely matching the non-adaptive results of 0.49 and 0.01 in Table 1. It is worth noting that the increased similarity between poisoned and correct texts in the adaptive setting leads to generally lower PAD values compared to Table 1, evidenced by a PAD of 1.88 for Vicuna-7B on MS-MARCO versus 1.91 in the non-adaptive setting.

### C.7 Impact of Poisoned Text Proportions in Retrieved Documents

To further evaluate the effect of varying numbers of poisoned texts in the retrieved results, we fix the number of retrieved texts to  $Top-K = 8$  and gradually increase the proportion of poisoned texts among the retrieved documents from 0% to 100%, where the number of poisoned texts is obtained by multiplying  $Top-K$  by the poisoning ratio and rounding down to the nearest integer. As shown in Table 13, PRA-RAG demonstrates strong robustness across varying poisoning ratios. At a poisoning rate of 20%, the system satisfies our robustness assumptions and achieves optimal defense performance with an ASR close to 0. When the poisoning rate increases to 30%–40%, exceeding the scope of our theoretical assumptions, the method remains empirically effective. As the poisoning rate continues to rise, the system avoids complete collapse even under an extreme poisoning scenario of 90%, evidenced by the Mistral-7B model on the NQ dataset where ACC decreases from 0.62 to 0.48 and ASR rises to 0.18. This phenomenon may stem from the ability of the LLM’s parametric knowledge to resolve conflicting contexts. Furthermore, the PAD metric exhibits a trend of initially increasing and then decreasing throughout this pro-

Table 13: The table shows the performance of PRA-RAG with varying ratios of poisoned texts ( $Top-K = 8$ , subset size = 3).

Models	Poisoning Rate	NQ			MS-MARCO			HotpotQA		
		ACC	ASR	PAD	ACC	ASR	PAD	ACC	ASR	PAD
Mistral-7B	0%	0.62	0.00	0.79	0.62	0.00	0.88	0.34	0.00	0.79
	10%	0.62	0.00	0.79	0.62	0.00	0.88	0.34	0.00	0.79
	20%	0.62	0.01	1.23	0.63	0.01	1.37	0.37	0.03	1.33
	30%	0.64	0.06	1.60	0.68	0.02	1.60	0.34	0.10	1.61
	40%	0.60	0.08	1.73	0.82	0.04	1.67	0.36	0.12	1.64
	50%	0.56	0.12	1.69	0.72	0.02	1.66	0.34	0.16	1.58
	60%	0.56	0.12	1.69	0.72	0.02	1.66	0.34	0.16	1.58
	70%	0.58	0.14	1.65	0.68	0.02	1.58	0.42	0.10	1.53
	80%	0.62	0.24	1.44	0.60	0.08	1.50	0.46	0.12	1.32
	90%	0.48	0.18	1.11	0.58	0.12	1.26	0.34	0.16	1.01
	100%	0.54	0.14	0.85	0.66	0.08	0.87	0.40	0.16	0.71
Vicuna-7B	0%	0.50	0.00	1.16	0.68	0.00	1.24	0.48	0.00	1.14
	10%	0.50	0.00	1.16	0.68	0.00	1.24	0.48	0.00	1.14
	20%	0.49	0.01	1.74	0.68	0.01	1.91	0.49	0.00	1.87
	30%	0.46	0.08	2.16	0.64	0.04	2.18	0.44	0.12	2.16
	40%	0.44	0.04	2.29	0.70	0.04	2.21	0.50	0.20	2.21
	50%	0.46	0.04	2.25	0.66	0.06	2.18	0.50	0.22	2.12
	60%	0.46	0.04	2.25	0.66	0.06	2.18	0.50	0.22	2.12
	70%	0.46	0.06	2.25	0.76	0.04	2.07	0.48	0.20	2.09
	80%	0.56	0.06	1.97	0.72	0.08	1.95	0.40	0.22	1.76
	90%	0.58	0.12	1.53	0.76	0.02	1.65	0.48	0.20	1.33
	100%	0.52	0.12	1.16	0.74	0.04	1.13	0.46	0.20	0.93

cess. Specifically, at medium poisoning rates (40%–60%), the adversarial interaction between clean and poisoned distributions causes the certified radius (PAD) to peak, while the significant drop in PAD at 100% poisoning indicates that poisoned samples are forced to form a tighter clustering structure to satisfy retrieval relevance constraints.

### C.8 Evaluation of the Impact of Monte Carlo Sampling

As the number of retrieved documents ( $Top-K$ ) and the subset size ( $n$ ) increase, the number of candidate combinations grows rapidly, incurring significant computational overhead. To reduce costs and enhance the practicality of our method, we introduce Monte Carlo sampling when the combination size reaches a certain level. Based on the analysis in Section A.3, we set the sampling count to  $m = 200$ , which also serves as the threshold for triggering the sampling process. Specifically, when the number of candidate combinations does not exceed  $m$ , we employ the original full enumeration strategy; when the number exceeds  $m$ , we utilize Monte Carlo sampling to approximate the combinatorial space. This strategy effectively controls computational overhead while minimizing the impact on performance.

Tables 14 and 15 present the performance of our method using Monte Carlo sampling approximation across various settings of retrieved document counts ( $Top-K$ ) and subset sizes ( $n$ ), specifically when the number of candidate combinations ex-

Table 14: Performance of PRA-RAG using Monte Carlo sampling under varying Top- $K$  (subset size = 3, poisoning rate = 10%).

Models	Top- $K$	NQ			MS-MARCO			HotpotQA		
		ACC	ASR	PAD	ACC	ASR	PAD	ACC	ASR	PAD
Mistral-7B	8	0.70	0.00	1.28	0.62	0.00	1.31	0.34	0.12	1.32
	12	0.52	0.00	1.16	0.66	0.00	1.21	0.44	0.02	1.18
	16	0.52	0.00	1.12	0.68	0.00	1.20	0.40	0.02	1.15
Vicuna-7B	8	0.44	0.00	1.78	0.66	0.02	1.84	0.42	0.04	1.86
	12	0.51	0.00	1.65	0.76	0.00	1.72	0.52	0.01	1.74
	16	0.42	0.00	1.64	0.70	0.00	1.70	0.48	0.00	1.67

Table 15: Performance of PRA-RAG using Monte Carlo sampling with Top- $K = 12$  retrieval under a poisoning rate of 10%, evaluated across different subset sizes.

Models	Subset Size	NQ			MS-MARCO			HotpotQA		
		ACC	ASR	PAD	ACC	ASR	PAD	ACC	ASR	PAD
Mistral-7B	3	0.52	0.00	1.16	0.66	0.00	1.21	0.44	0.02	1.18
	4	0.56	0.00	1.02	0.66	0.00	1.07	0.44	0.02	1.07
	5	0.72	0.00	0.90	0.82	0.00	0.97	0.46	0.01	0.96
Vicuna-7B	3	0.51	0.00	1.65	0.76	0.00	1.72	0.52	0.01	1.74
	4	0.48	0.00	1.44	0.66	0.02	1.51	0.47	0.02	1.55
	5	0.47	0.00	1.30	0.60	0.00	1.37	0.45	0.02	1.36

ceeds the threshold ( $m = 200$ ). By comparing these results with those obtained via full enumeration in Tables 6 and 8, it is evident that the introduction of Monte Carlo sampling causes no significant degradation in performance. Furthermore, the overall trends remain consistent, indicating that the impact of poisoned texts is further attenuated as Top- $K$  and  $n$  increase, while PAD exhibits a continuous downward trend.

Tables 16 and 17 present the computational overhead of the algorithm with Monte Carlo sampling under different Top- $K$  and  $n$  settings. In Table 16, with Top- $K = 8$  and  $n = 3$ , the number of candidate combinations is 56, which remains below the threshold ( $m = 200$ ), implying that sampling is not triggered and the computational cost remains consistent with the results in Table 10. In other scenarios, compared to the full enumeration results in Tables 10 and 11, the sampling strategy

Table 16: Computational overhead of PRA-RAG using Monte Carlo sampling under varying Top- $K$  (subset size  $n = 3$ , poisoning rate = 20%).

Models	Top- $K$	NQ	MS-MARCO	HotpotQA
Mistral-7B	8	5.98s	6.20s	6.17s
	12	11.78s	11.71s	11.97s
	16	11.09s	11.69s	12.50s
Vicuna-7B	8	3.37s	2.91s	3.02s
	12	6.67s	7.55s	6.42s
	16	9.23s	7.84s	6.46s

Table 17: Computational overhead of PRA-RAG using Monte Carlo sampling with Top- $K = 12$  under a poisoning rate of 10%.

Model	Subset Size	NQ	MS-MARCO	HotpotQA
Mistral-7B	3	12.26s	10.51s	11.93s
	4	11.96s	15.07s	12.16s
	5	11.95s	15.09s	12.17s
Vicuna-7B	3	6.87s	9.03s	6.79s
	4	7.03s	7.70s	6.70s
	5	7.17s	8.09s	7.28s

significantly reduces computational overhead. For instance, on the Vicuna-7B model using the NQ dataset, the cost drops from 64.93s to 9.23s when Top- $K = 16$  and  $n = 3$ , and from 116.09s to 7.17s when Top- $K = 12$  and  $n = 5$ .

1244  
1245  
1246  
1247  
1248