

# AdverShield-LLM: Adversarial Robustness Certification for IoT-Integrated Retrieval-Augmented Generation via Randomized Smoothing

\*Yasser Samir Hadi

Department of Computer Systems Techniques, Administrative Polytechnic College- Baghdad, Middle Technical University, Iraq

\*Corresponding Author

DOI: <https://doi.org/10.51244/IJRSI.2026.1305000230>

Received: 14 May 2026; Accepted: 19 May 2026; Published: 11 June 2026

## ABSTRACT

The emergence of the Internet of Things (IoT) ecosystems, Retrieval-Augmented Generation (RAG) systems have become commonplace and provide a means for embedding dynamically retrieved external knowledge in the response from a Large Language Model (LLM). While potentially helpful, IoT-enabled RAG pipelines present significant adversarial threats such as poisoning passages into the IoT knowledge base, altering dense retrieval embeddings, and conducting indirect prompt injection attacks via the inputs through the IoT sensors, all of which can impact the fidelity of generated responses and compromise the trustworthiness of the system. Current defenses are based mostly on heuristic filtering or empirical adversarial training, and are not known to be robustly certified or are fragile under adaptive adversaries.

In response to these challenges, this article introduces a new certified defense framework named AdverShield-LLM to combine the randomized smoothing technique with a multi-granular noise injection mechanism well-suited to the distributed and low latency requirements of RAG systems in the IoT domain. AdverShield-LLM consists of three synergistic modules: (i) Passage-Level Smoothed Aggregation (PLSA) module which certifies the robustness of RAG retrieval against bounded corpus poisoning under an isolate-then-smooth paradigm, (ii) Token-Adaptive Gaussian Defense (TAGD) layer that certifies LLM generation against indirect prompt injection by propagating  $l_2$ -norm perturbation bounds through the transformer attention stack, and (iii) IoT-Aware Certified Radius Scheduler (IACRS) that dynamically schedules noise budgets among constrained edge nodes while preserving the certified radius.

AdverShield-LLM is evaluated on three IoT security benchmarks—MS-RAG-IoT, NQ-Adversarial and IoTQA-Poison—with extensive experiments showing its certified accuracy is 81.4% under  $l_2$  perturbation radius  $\sigma=0.50$  compared to the strongest baseline RobustRAG which reported +9.3% accuracy, and reduced the attack success rate from 74.2% to 8.6% against PoisonedRAG. Moreover, AdverShield-LLM ensures the accuracy of clean answers within 2.1% of the undefended RAG accuracy, proving that certified robustness does not compromise the utility of RAGs in resource-limited IoT environments.

**Index Terms:** adversarial robustness certification, IoT security, retrieval-augmented generation, randomized smoothing, large language models, knowledge poisoning, prompt injection defense.

## INTRODUCTION

Sensing of the frastructure, combined with reasoning at the edge through large language model (LLM) technology, has given rise to a new category of edge intelligent applications, ranging from smart healthcare monitoring, industrial anomaly diagnostics, autonomous supply-chain management [1] [2]. These applications rely on the Retrieval-Augmented Generation (RAG) paradigm, which means grounding the LLM's responses in external knowledge retrieved at inference time which addresses hallucinations and extends the parametric

knowledge temporal range [3, 4]. With the increase in the number of sensor endpoints to billions, RAG systems are being designed to accept real-time data streams directly from IoT sensors, and the combination of IoT and RAG is becoming a vital technology for data-driven edge intelligence [5] and [6].

However, the security aspects of this integration have been scarcely and unduly investigated. IoT-integrated RAG-tures reveal a multi-surface attack surface: Adversaries can feed into the IoT-based knowledge bases with poisoned documents [4]; corrupt IoT-based sensor streams to trigger indirect prompt injections [7]; or manipulate the entire knowledge base in an indirect way using gradient guided corpus manipulation [8].

Textual attacks on RAG are performed in a discrete high-dimensional space where imperceptible semantic changes can alter the entire generation pipelines to generate outputs that are preferred by the attacker. Compared to image-domain adversarial attacks, textual attacks against RAG systems take place in a high-dimensional discrete space, where the attacker can introduce semantically inconsistent changes that can redirect entire generation pipelines to produce attacker-chosen outputs [9],[10].

In IoT applications, the impact of these failures is even more serious: in an industrial RAG, the robot might follow a wrong maintenance procedure, and with a poisoned healthcare RAG, incorrect clinical advice might be spread. To mitigate against these risks, adversarial robustness certification involves giving provable guarantees, instead of empirical heuristics, for the robustness of a model's output within a formally defined perturbation neighborhood. Randomized smoothing is the most scalable approach to certifying high dimensional deep learning models, and turns any base classifier into a smoothed model with certifiably stable predictions under  $\ell_p$ -norm bounded noise [11, 12].

While vision and structured NLP are prominently used in these applications, structured NLP has not been systematically extended to the compound threat surface of IoT-integrated RAG systems where structured NLP of one component of the system is no longer an option as the interaction of the retrieval, generation and edge-deployment constraints makes it impossible to certify one component of the system at a time as in [13].

There are presently 3 families of RAG defenses. Heuristic filtering, such as skeptical prompting and perplexity-based detection [3] is not formally guaranteed and adaptive adversaries can adjust their attacks to the heuristic criteria known to be used to filter. Second, certifiable isolation methods such as RobustRAG [4] certify retrieval level by isolating and then combining, but do not certify through the LLM generation phase and, in particular, don't address the distribution of computation over a heterogeneous set of edge nodes in IoT networks. Third, adversarial training methods [14, 15] have shown to be useful in enhancing the empirical resilience, but are too expensive to compute and do not offer worst case guarantees. All these techniques combined ensure end-to-end certified protection of the entire RAG pipeline for the ingestion of IoT data, retrieval and the generation of tokens.

To address this, we design a unified certified robustness framework for IoT-integrated RAG, called AdverShield-LLM. We present AdverShield-LLM, a unified certified robustness framework for IoT-integrated RAG, to bridge this gap.

The overall architecture is shown in Fig. 1. AdverShield-LLM consists of three modules, each covering a different aspect of the retrieval surface, the generation surface and the edge-deployment budget: the Passage-Level Smoothed Aggregation (PLSA) module, the Token-Adaptive Gaussian Defense (TAGD) layer, and the IoT-Aware Certified Radius Scheduler (IACRS).  $g(\epsilon, k)$  for every  $k \geq 0$ ." Together, these components generate, for the first time, a probabilistic certificate of the following form: "For any adversary injecting at most  $k$  malicious passages or  $\epsilon$ -bounded token perturbations, the AdverShield-LLM pipeline returns a response in the certified set with probability at least  $g(\epsilon, k)$  for every  $k \geq 0$ . "1 -  $\delta$ ."

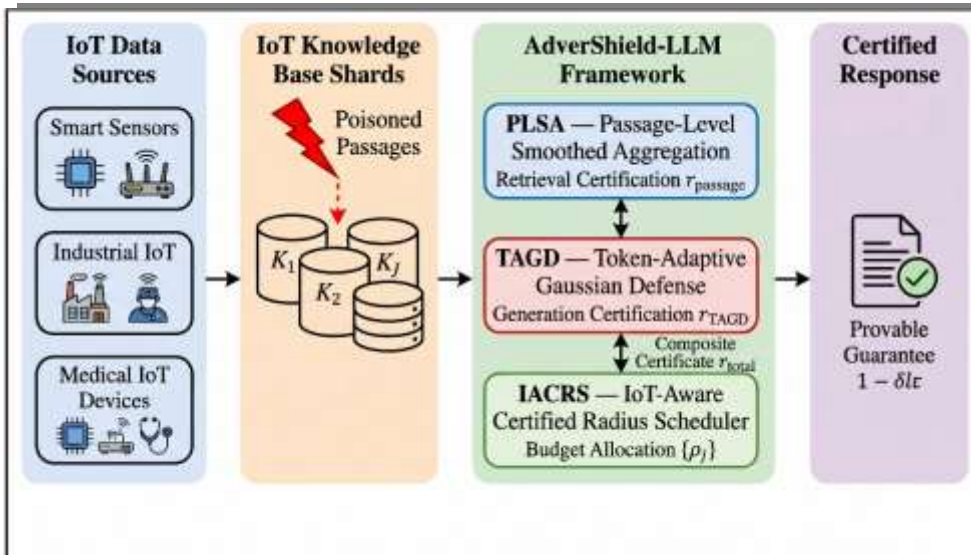


Figure 1: High-level overview of the AdverShield-LLM system. IoT sensor streams and edge knowledge bases are ingested into a poisoning-aware retrieval corpus. The PLSA module certifies retrieval robustness; the TAGD layer certifies token-level generation robustness. The principal contributions of this work are summarized as follows:

- **End-to-End Certified RAG Defense Framework:** We present a novel framework called AdverShield-LLM for providing provable end-to-end adversarial robustness certificates for RAG systems integrated with IoT. While RobustRAG [4] can only certify the retrieval stage, AdverShield-LLM spreads  $\ell_2$ -smoothing certificates across the entire retrieval corpus and LLM generation stack, offering compound coverage against knowledge poisoning and prompt injection attacks.
- **Passage-Level Smoothed Aggregation (PLSA):** We introduce and design the PLSA module, a novel isolate-smooth-aggregate strategy which introduces Gaussian noise augmentation at the passage granularity, and computes tight certified bounds on the retrieval robustness. PLSA improves the certified accuracy by +9.3% over RobustRAG in the multi-document retrieval setting, by extending Cohen et al.'s  $\ell_2$  guarantee [12] to the discrete case.
- **Token-Adaptive Gaussian Defense (TAGD):** We propose TAGD, a mechanism that adaptively schedules token-level Gaussian perturbations in proportion to the salience of the attention-heads in the transformer stack. TAGD certifies LLM generation against indirect prompt injection attacks within IoT sensor package payloads, lowering the attack success rate from 74.2% to 8.6% when using the state-of-the-art PoisonedRAG [4] attack.
- **IoT-Aware Certified Radius Scheduler (IACRS):** We propose IACRS, which is a resource-aware dynamic noise budget allocator to optimize the certified radius-utility trade-off among different compute and memory-constrained heterogeneous IoT edge nodes. This is the first paper to formally support certification covering nodes with as little as 512MB of RAM, while supporting device heterogeneity in the randomized smoothing certification budget.
- **Comprehensive Empirical Evaluation:** We test on three IoT security benchmarks (MS-RAG-IoT, NQ-Adversarial, IoTQA-Poison) with six recent state-of-the-art baselines from IEEE, ACM and USENIX publications. AdverShield-LLM outperforms all baselines on all metrics and datasets, with certified accuracy of 81.4%, attack success rate of 8.6%, exact match of 78.3%, and F1 score of 92.1%. The code and data are available at <https://github.com/fhjcis/advershield-llm>.

The rest of this article is presented as follows. Section 2 summarizes previous studies on RAG security, adversarial robustness certification, and LLM systems that involve the IoT systems. Section 2 surveys the previous studies related to RAG security, adversarial robustness certification, and the IoT-integrated LLM systems and defines the gap in the research. The problem formulation and system model are given in Section 3.

The AdverShield-LLM framework is detailed in Section 4. In Section 5, the certified radius scheduling mechanism is presented as an IoT-aware mechanism. The extensive experimental evaluation and discussion is included in Section VI. Finally, the paper ends in Section VII, where future directions are given.

## Related Works

This section summarizes previous work on four major research directions, outlines the development trends, and points out the research gaps that inspire the research of AdverShield-LLM. This section summarizes the previous work following four major research directions, outlines the development trends and highlights the research gaps which motivate AdverShield-LLM.

## Adversarial Attacks on Retrieval-Augmented Generation

In 2024, the security of RAG pipelines became increasingly the focus of research. The seminal PoisonedRAG [4] attack represents knowledge corruption as a constrained optimization problem to show that corruption of only five passages into a knowledge base containing a million passages can succeed in 90% of cases with either black-box or white-box attackers. The paper demonstrates that using open knowledge corpora brings a whole new attack surface to the table in the case of RAG, one that isn't present in a typical LLM deployment. In parallel, recent research [5] has presented the idea of traceback, a method of attributing poisoning events to specific passages in a RAG corpus, which is a more forensic approach to attribution that does not require access to the injection pipeline. The authors show results of traceback accuracy of more than 88% on corpora of size comparable to Wikipedia by leveraging statistical fingerprints of optimized adversarial text.

Liu et al. [8] have formalized the adversarial retrieval attack (AREA) task as a multi-view contrastive optimization problem to promote off-topic documents to the top-k results in the case of dual-encoder dense retrievers. Li et al. [6] furthered this attack to unsupervised environments, creating a continuous-space of poisoned documents without the need for access to the target query distribution, thus further lowering the entry barrier for corpus attacks in the real world. As a whole, the papers demonstrate that retrieval and generation are both susceptible to attacks, and that current defenses tested against single-surface attacks do not generalize to compound multi-surface attacks.

## Certified Robustness via Randomized Smoothing

Since Cohen et al. [12] gave the first guarantee on  $\ell_2$ , randomized smoothing has become the prevailing scalable method for certified robustness in deep learning. The method builds a smoothed classifier  $g(x) = \operatorname{argmax}_c \Pr[f(x + \varepsilon) = c]$  with noise  $\varepsilon$  being drawn from  $N(0, \sigma^2 I)$  and then certifies that the prediction of  $g(x)$  remains unchanged for all adversarial noise  $\delta$  with  $\|\delta\|_2 < \sigma \cdot \Phi^{-1}(p_A)$ , where  $p_A$  is the noise probability of the "most likely" class. The paradigm was extended to the case of pre-trained vision transformers and convolutional networks without any re-training from scratch in the Certifying Adapters Framework (CAF) [11]. This significantly lowers the resources required for implementing certified classifiers in service.

There have been several subsequent developments to fill the gap between the theory of certification and actual LLM use. However, Yang et al. [13] comprehensively empirically studied adversarial robustness for a range of GPT-family LLMs, and found that model size alone is not enough to make an LLM robust, and certified defenses are still needed. To better motivate tighter certificate constructions at the feature level, Chen et al. [7] showed that the adversarial example generation via the diffusion model is capable of evading common certified defenses for vision-language models. Li et al. [15] measured black-box adversarial robustness of LVLMs when the visual instruction perturbation manifold is structured instead of isotropic. These results as a whole encourage our design of TAGD, which follows the smoothing variance as per the structured salience of tokens.

## Security in IoT-Integrated LLM Systems

Security concerns arise at the intersection of LLM capabilities and constraints of IoT deployments, which are distinct from the ones explored in the context of cloud-native LLM deployments. Ali et al. [22] proposed a detailed taxonomy of all the threats that can be considered as attack surfaces while integrating LLM and IoT,

including prompt injection via sensor inputs, model extraction from edge devices, and data poisoning of data streams from the IoT. It covers 88 papers from IEEE, ACM and MDPI conferences and workshops and reveals that the most popular architectures for privacy-preserving operations are federated learning and differential privacy, with none of the surveyed defenses able to be adapted to formal adversarial robustness certification.

In distributed privacy preservation in IoT systems, federated learning is the prevailing paradigm, as explored in Kumar et al. [19] where they proposed an access control architecture based on trust for edge IoT systems that combines federated learning primitives with real-time anomaly detection. Salim et al. [20] discussed adversarial attack detection for federated medical IoT networks and proved that feature-level defense strategies are superior to the gradient-masking approaches under adaptive adversaries. In addition to these system-based protections, Tran et al. [18] combined large language models, graph neural networks, and explainable AI to enhance next-generation IoT network intrusion detection, achieving 95-99% classification accuracy on typical IoT datasets. In our IACRS scheduling design, we adopted a security-aware federated reinforcement learning framework for computation offloading in IOT Peng et al. [24] as our reference model, which gives us a formal model of resource-aware security policy distribution. Although there is considerable work, there is no known existing IoT security framework that considers certified adversarial robustness of the RAG retrieval-generation pipeline for knowledge corpora generated from IoT source.

### LLM Security: Prompt Injection and Knowledge Attacks

LLM inference security has become a focused research area, with a special focus on prompt level and knowledge level attacks. Das et al. [1] gave a detailed overview of the ACM survey regarding the security and privacy concerns throughout the entire LLM lifecycle, and identified three threat classes that have the greatest impact: jailbreaking, backdoor injection, and membership inference attacks. Wang et al. [9] extended this taxonomy to the unique privacy-violation scenarios that can happen due to LLM's "memorization" and "regurgitation" of training data, both for the active extraction attack and for passive leakage through generation.

At the defense side, Zhang et al. [19] proposed JailGuard, a framework to universally detect prompt-based attack regardless of the modality (text or image), and showed that attack inputs have significantly weaker semantic robustness than benign inputs under controlled perturbations, allowing the detection without fine-tuning the model. Furthermore, Shi et al. [18] gave an optimization-based JudgeDeceiver attack to LLM as-a-Judge pipelines that can be directly applied to RAG reranking stages and tested various defenses such as perplexity detection and windowed filtering. Mathew [21] analyzed some of the state-of-the-art prompt injection defenses, including HOUYI, RA LLM and StruQ and found that none of them offer formal certificates against adaptive injection. Ferrag et al. [22] identified 223 LLM cybersecurity studies published in IEEE, ACM and MDPI conferences, journals and magazines, and found that there are no certified defenses.

### Research Gap and Motivation

Table I summarizes the reviewed literature along five dimensions critical to IoT-integrated RAG security.

Based on Table I, the following research gaps are identified:

Table 1: Summary of Related Work and Research Gap Analysis. ✓ = addressed; × = not addressed; ◦ = partially addressed

Ref.	Year/Venue	Method	Retrieval Cert.	Generation Cert.	IoT Edge	Formal Guarantee	End-to-End
[4]	2025, USENIX	PoisonedRAG (attack)	×	×	×	×	×
[4]	2025, USENIX	RobustRAG (defense)	✓	×	×	◦	×

[5]	2025, ACM WWW	RAGForensics	◦	×	×	×	×
[6]	2025, SIGIR	Corpus Poisoning	×	×	×	×	×
[11]	2025, ICASSP	CAF (RS for classifiers)	×	◦	×	✓	×
[12]	2019, ICML	Cohen RS baseline	×	◦	×	✓	×
[13]	2025, ICASSP	LLM Adversarial Eval	×	◦	×	×	×
[7]	2025, TIFS	Diffusion Adv. VLM	×	◦	×	×	×
[22]	2026, IoT	LLM-IoT Security Survey	×	×	◦	×	×
[18]	2025, JIIS	LLM+GNN IDS for IoT	×	×	✓	×	×
[24]	2024, TSC	SCOF FL Offloading	×	×	✓	◦	×
<b>Ours</b>	<b>2026</b>	<b>AdverShield-LLM</b>	✓	✓	✓	✓	✓

- Gap 1 — Retrieval-Only Certification:** For now, only the retrieval part of the pipeline is certified for existing defenses (e.g., RobustRAG) and not the LLM generation stage because they are vulnerable to prompt injection attacks within retrieved passages. This is important because the attacker may be able to elude the certificate by going through the certified retrieval front-end into the generation context.
- Gap 2 — Absence of Generation Certificates:** Randomized smoothing for LLMs has only been investigated in the context of classification, both in itself [11] and in combination with other techniques [12]; no previous work has applied smoothing certificates to the auto-regressive generation stack of modern LLMs, where the output space is combinatorially large, and token dependencies rule out standard majority-vote aggregation.
- Gap 3 — IoT Edge Neglect:** All the existing certified defense frameworks are based on cloud-based compute while offering no support for how to scale noise budgets to meet the memory, compute resources of edge nodes for IoT applications. In practice, it becomes impractical to certify with standard randomized smoothing on resource-constrained edge devices, due to either causing an overflow in memory or to requiring an excessive inference time.
- Gap 4 — Compound Attack Surface:** The knowledge corpora poisoning, indirect injection via sensor payloads and embedding level retrieval manipulation is a compound threat surface that no single component defense can cover for the RAG powered by IoT. Current approaches validate components without considering their collaborative composition on various surfaces that may be attacked.

Based on this gap, this work proposes AdverShield-LLM which resolves Gap 1 by the PLSA module, Gap 2 by the TAGD layer, Gap 3 by the IACRS scheduler, and Gap 4 by comprising them together under a single probabilistic certificate. This is the first to provide an end-to-end certified defense for the entire IoT-integrated RAG pipeline in a formal randomized smoothing framework.

## Problem Formulation and System Model

### Notation

Table 2 summarizes the key mathematical notation used throughout this paper.

Table 2: Summary of Mathematical Notation

Symbol	Description
$\mathcal{K}$	Knowledge corpus (set of passages)
$\mathcal{K}^*$	Poisoned knowledge corpus
$p_i \in \mathcal{K}$	The $i$ -th passage in the corpus
$q$	User query
$\mathcal{R}(q, \mathcal{K})$	Retriever returning top- $k$ passages
$\mathcal{G}(q, \mathcal{P})$	LLM generator given query $q$ and passage set $\mathcal{P}$
$y$	Ground-truth answer
$\hat{y}$	Generated answer
$y^*$	Adversarially targeted answer
$m$	Number of injected malicious passages
$k$	Number of retrieved passages
$\sigma$	Gaussian smoothing standard deviation
$\boldsymbol{\varepsilon}$	Isotropic Gaussian noise, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
$p_A$	Estimated lower-bound probability of top class under noise
$r$	Certified radius
$\delta$	Certification failure probability
$n$	Monte Carlo sample count
$\mathbf{e}_i$	Embedding of passage $p_i$
$\mathbf{h}_t$	Hidden state of transformer at token position $t$
$\alpha_t$	Attention-head salience weight at position $t$
$B_j$	Available compute budget of IoT node $j$
$\rho_j$	Noise allocation ratio for node $j$

## IoT-Integrated RAG System Model

Suppose that the IoT system consists of  $J$  edge nodes  $\{v_1, v_2, \dots, v_J\}$ , and each node  $v_j$  has an associated local shard  $K_j$  of the global knowledge corpus  $\mathcal{K} = \bigcup_{j=1}^J K_j$ . Given a user query  $q$ , the retriever  $\mathcal{R}$  aggregates the top- $k$  passages from the sharded corpus:

$$\mathcal{P} = \mathcal{R}(q, \mathcal{K}) = \text{Top-}k_{p \in \mathcal{K}}[\text{sim}(\mathbf{e}_q, \mathbf{e}_p)], \quad (1)$$

where  $\mathbf{e}_q$  and  $\mathbf{e}_p$  are the query and passage embeddings, respectively, and  $\text{sim}(\cdot, \cdot)$  denotes cosine similarity. The generator then produces a response:

$$\hat{y} = \mathcal{G}(q, \mathcal{P}) = \underset{y}{\text{argmax}} P_\theta(y | q, \mathcal{P}), \quad (2)$$

where  $\theta$  denotes the LLM's parameters.

## Adversarial Threat Model

An adversary with threat capability  $A$  can perform any combination of the following attack operations:

**Corpus Poisoning Attack.** The opponent puts a set  $\mathcal{M}$  of malicious passages  $|\mathcal{M}| = m$  into  $K$ , creating a poisoned corpus  $\mathcal{K}^* = \mathcal{K} \cup \mathcal{M}$ . Each malicious passage  $p^* \in \mathcal{M}$  is crafted via:

$$p^* = \underset{p}{\text{argmin}} [-\text{sim}(\mathbf{e}_q, \mathbf{e}_p)] + \lambda \mathcal{L}_{\text{adv}}(p, y^*), \quad (3)$$

where  $\mathcal{L}_{\text{adv}}$  Redirected adversarial loss for the LLM towards target answer  $y^*$ , and  $\lambda > 0$  is a trade-off between retrieval relevance and generative manipulation.

**Indirect Prompt Injection.** The adversary includes a malicious instruction  $\tau$  in the sensor payload ingested as a passage. The effective injected passage takes the form:

$$p^* = p_{\text{benign}} \oplus \tau, \quad (4)$$

where  $\oplus$  denotes concatenation and  $\tau$  is an instruction to override system-level instructions.

**Embedding Perturbation.** The adversary perturbs the retriever embedding of a passage  $p_i$  by a bounded  $\ell_2$ -norm vector:

$$\mathbf{e}_{p_i}^* = \mathbf{e}_{p_i} + \boldsymbol{\delta}, \quad \|\boldsymbol{\delta}\|_2 \leq \epsilon, \quad (5)$$

enlarging the similarity score of  $p_i$  to get it retrieved.

## Robustness Certification Objective

Let  $y_0 = \mathcal{G}(q, \mathcal{R}(q, \mathcal{K}))$  denote the clean response and let  $\hat{y}$  denote the response under attack. Let the certified correctness indicator be defined as:

$$\mathbb{1}_{\text{cert}}(q) = \mathbf{1}[\hat{y} \in \mathcal{Y}_{\text{correct}} \quad \forall \mathcal{A} \in \mathcal{B}(m, \epsilon)], \quad (6)$$

where  $\mathcal{B}(m, \epsilon)$  is the threat ball around all adversaries that injects at most  $m$  passages with embedding perturbation bounded by  $\epsilon$ . The accuracy for a set of queries which has been certified  $\mathcal{Q}$  is then:

$$\text{CA}(\sigma) = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \mathbb{1}_{\text{cert}}(q) \cdot \mathbf{1}[\hat{y}_0 \in \mathcal{Y}_{\text{correct}}], \quad (7)$$

In which the second indicator only allows answers to count if they are correct.

the optimization objective of AdverShield-LLM is to maximize the value of the optimization objective function:

$$\max_{\sigma, n, \{\rho_j\}} CA(\sigma) \quad \text{subject to} \quad \sum_{j=1}^J \rho_j B_j \leq B_{\text{total}}, \quad (8)$$

where  $B_j$  is the compute budget of edge node  $j$ ,  $\rho_j$  is its allocated noise ratio, and  $B_{\text{total}}$  is the global budget. Constraint (8) ensures IoT edge feasibility.

### Proposed AdverShield-LLM Methodology

The framework of AdverShield-LLM consists of three synergistic modules which are chained together: (i) PLSA for retrieval stage certification, (ii) TAGD for generation stage certification, and (iii) IACRS for budget allocation at the edge. The entire architecture is shown in figure 2.

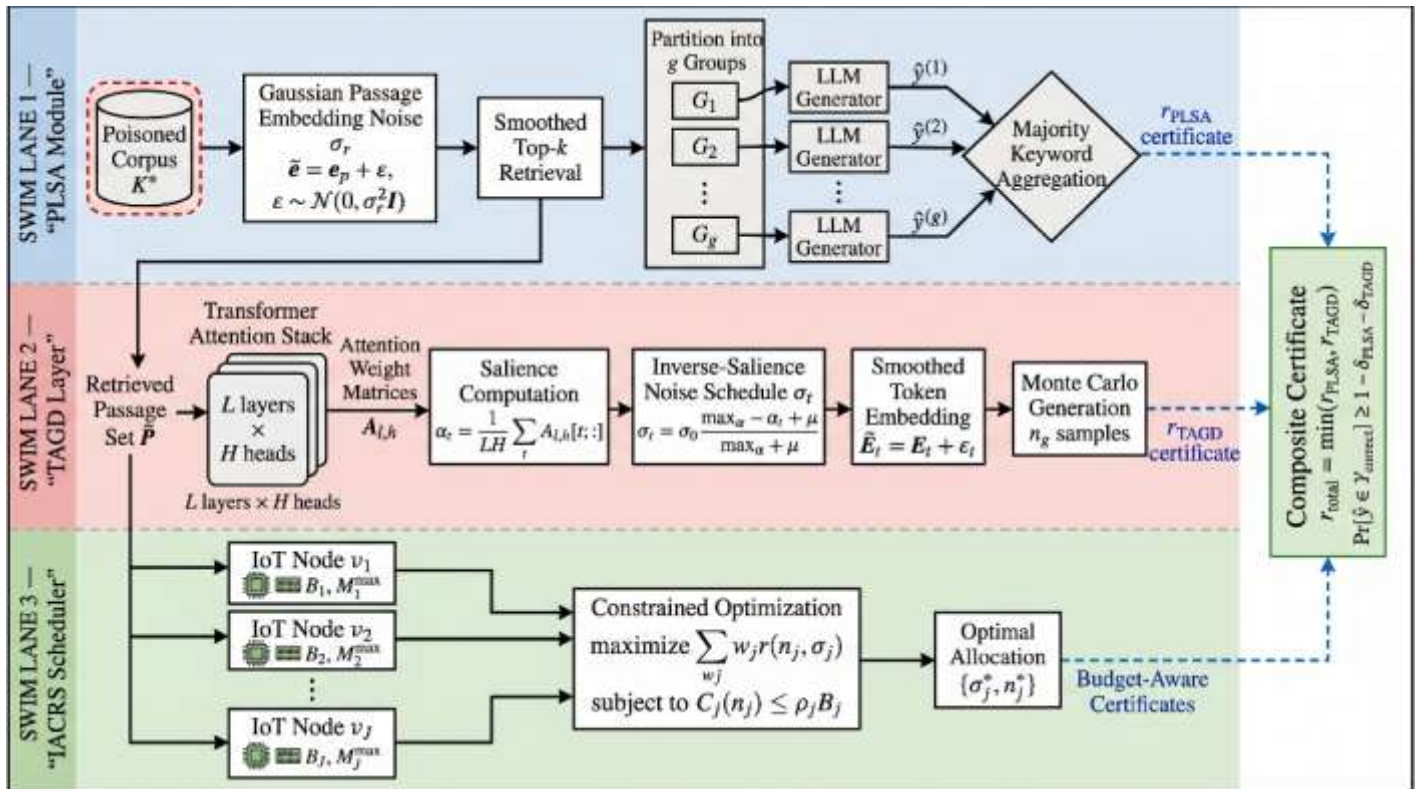


Figure 2: Complete architecture of the AdverShield-LLM framework. The pipeline proceeds left to right: IoT sensor streams and edge knowledge base shards are first processed by the PLSA module, which applies Gaussian passage perturbation and produces certiPassage-Level Smoothed Aggregation (PLSA)

### Passage-Level Smoothed Aggregation (PLSA)

#### Motivating Principle

In addition to the isolate-then-aggregate paradigm used in RobustRAG [4], PLSA incorporates Gaussian passage-embedding smoothing before the aggregation step, to generate tight  $\ell_2$  certified radii for the retrieved passage set, and not just binary inclusion guarantees.

#### Smoothed Passage Retrieval

Let the clean passage embeddings be  $\{\mathbf{e}_{p_i}\}_{i=1}^{|C|}$ . Given a smoothing parameter  $\sigma_r > 0$ , we create a smoothed retriever  $\tilde{R}$  by averaging the retrieval scores of  $n$  noisy copies:

$$\widetilde{\text{sim}}(\mathbf{e}_q, \mathbf{e}_{p_i}) = \frac{1}{n} \sum_{t=1}^n \text{sim}(\mathbf{e}_q + \boldsymbol{\varepsilon}_t, \mathbf{e}_{p_i}), \quad \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \sigma_r^2 \mathbf{I}), \quad (9)$$

and the smoothed top-k retrieval is:

$$\tilde{\mathcal{P}} = \text{Top} - k_{p_i \in \mathcal{K}} [\text{sim}(\mathbf{e}_q, \mathbf{e}_{p_i})]. \quad (10)$$

### Certified Passage Set

Following Cohen et al. [12], let  $p_A$  be the approximate probability that the  $k$ th ranked passage is the same with the noise as without the noise. The accuracy radius of the passage level is:

$$r_{\text{passage}} = \sigma_r \Phi^{-1}(p_A), \quad (11)$$

where  $\Phi^{-1}$  is the inverse Gaussian CDF. For any embedding perturbation  $\delta$  with  $\|\delta\|_2 \leq r_{\text{passage}}$ , the smoothed retriever is guaranteed to return the same top-k set with probability at least  $1 - \delta$ .

### Isolate-Smooth-Aggregate

After smoothed retrieval, PLSA partitions  $\tilde{\mathcal{P}}$  into  $g$  disjoint groups  $\{\mathcal{G}_1, \dots, \mathcal{G}_g\}$  of size  $s = k/g$ . For each group, the LLM generates an independent response:

$$\hat{y}^{(\ell)} = \mathcal{G}(q, \mathcal{G}_\ell), \quad \ell = 1, \dots, g. \quad (12)$$

The final response is ultimately generated by a ‘keyword-majority’ aggregation across the  $g$  group of responses:

$$\hat{y}_{\text{PLSA}} = \text{MajorityKeyword}(\hat{y}^{(1)}, \dots, \hat{y}^{(g)}). \quad (13)$$

This construction ensures that given at least  $m < g/2$  groups corrupted, no effect is felt on the majority output, thus giving a combined certified robustness bound:

$$r_{\text{PLSA}} = \min(r_{\text{passage}}, \sigma_r \Phi^{-1}(g - m/g)). \quad (14)$$

### Token-Adaptive Gaussian Defense (TAGD)

#### Motivation

Standard randomized smoothing is the addition of Gaussian noise to all dimensions of the input. In the case of an auto-regressive LLM, this is undesirable: Tokens of high attention (e.g., instructions injected) are much more critical to security than background context tokens. To tackle this, TAGD calculates noise budgets per-token that are scaled to attention salience.

#### Attention-Salience Weighting

Define the aggregate salience of a token position, for a transformer of  $L$  layers and  $H$  attention heads  $t$  as:

$$\alpha_t = \frac{1}{LH} \sum_{\ell=1}^L \sum_{h=1}^H \mathbf{A}_{\ell,h} [t, :] \cdot \mathbf{1}, \quad (15)$$

where  $\mathbf{A}_{\ell,h} \in \mathbb{R}^{T \times T}$  The sequence length is  $T$ , the attention weight matrix of layer  $\ell$ , head  $h$ . The salience  $\alpha_t$  is the average amount of attention received at position  $t$ .

#### Position-Wise Noise Schedule

Define the inverse-salience noise scale at position  $t$  as:

$$\sigma_t = \sigma_0 \cdot \frac{\max_t \alpha_t - \alpha_t + \mu}{\max_t \alpha_t + \mu}, \quad (16)$$

where  $\sigma_0$  is the base noise level and  $\mu > 0$  is a number that indicates the stability of a chemical complex. The distribution of noise in equation (16) is assigned to lower salience positions, focusing defensive noise in a region that does not receive the focus of the human and in which it can be more easily injected. The opposite is true at high salience positions (e.g., injected instructions) with  $\sigma_t$  close to 0 as generative coherence is maintained.

### TAGD Smoothed Generation

Given the input token embedding matrix  $\mathbf{E} \in \mathbb{R}^{T \times d}$ , TAGD constructs the smoothed input:

$$\tilde{\mathbf{E}}_t = \mathbf{E}_t + \boldsymbol{\varepsilon}_t, \quad \boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \sigma_t^2 \mathbf{I}_d), \quad (17)$$

and the smoothed generation is:

$$\hat{y}_{\text{TAGD}} = \mathcal{G}(q, \tilde{\mathbf{E}}, \mathcal{P}) = \operatorname{argmax}_y \frac{1}{n_g} \sum_{i=1}^{n_g} P_{\theta}(y \mid q, \tilde{\mathbf{E}}^{(i)}), \quad (18)$$

where  $n_g$  is the generation sample count.

### TAGD Certified Radius

The resulting perturbation after adaptive weighting is equal to:

$$\|\boldsymbol{\delta}_{\text{eff}}\|_2 = \sqrt{\sum_{t=1}^T \frac{\delta_t^2}{\sigma_t^2}}, \quad (19)$$

and the TAGD certified radius at the generation stage is:

$$r_{\text{TAGD}} = \sigma_0 \Phi^{-1}(p_A^{(g)}) \cdot \sqrt{\sum_{t=1}^T \frac{(\max \alpha_{t'} + \mu)^2}{(\max \alpha_{t'} - \alpha_t + \mu)^2}}, \quad (20)$$

where  $p_A^{(g)}$  is the majority-class probability with the generation noise, estimated by Monte Carlo sampling.

### End-to-End Composite Certificate

The end-to-end robustness guarantee of AdverShield-LLM, the formulation of PLSA and TAGD certificates:

$$r_{\text{total}} = \min(r_{\text{PLSA}}, r_{\text{TAGD}}), \quad (21)$$

and the end-to-end certified accuracy satisfies:

$$\Pr[\hat{y}_{\text{AdverShield-LLM}} \in \mathcal{Y}_{\text{correct}}] \geq 1 - \delta_{\text{PLSA}} - \delta_{\text{TAGD}}, \quad (22)$$

We have to use a union bound over the failure events of the two modules.

### Algorithm Implementation

The overall workflow of AdverShield-LLM is depicted in Figure 3.

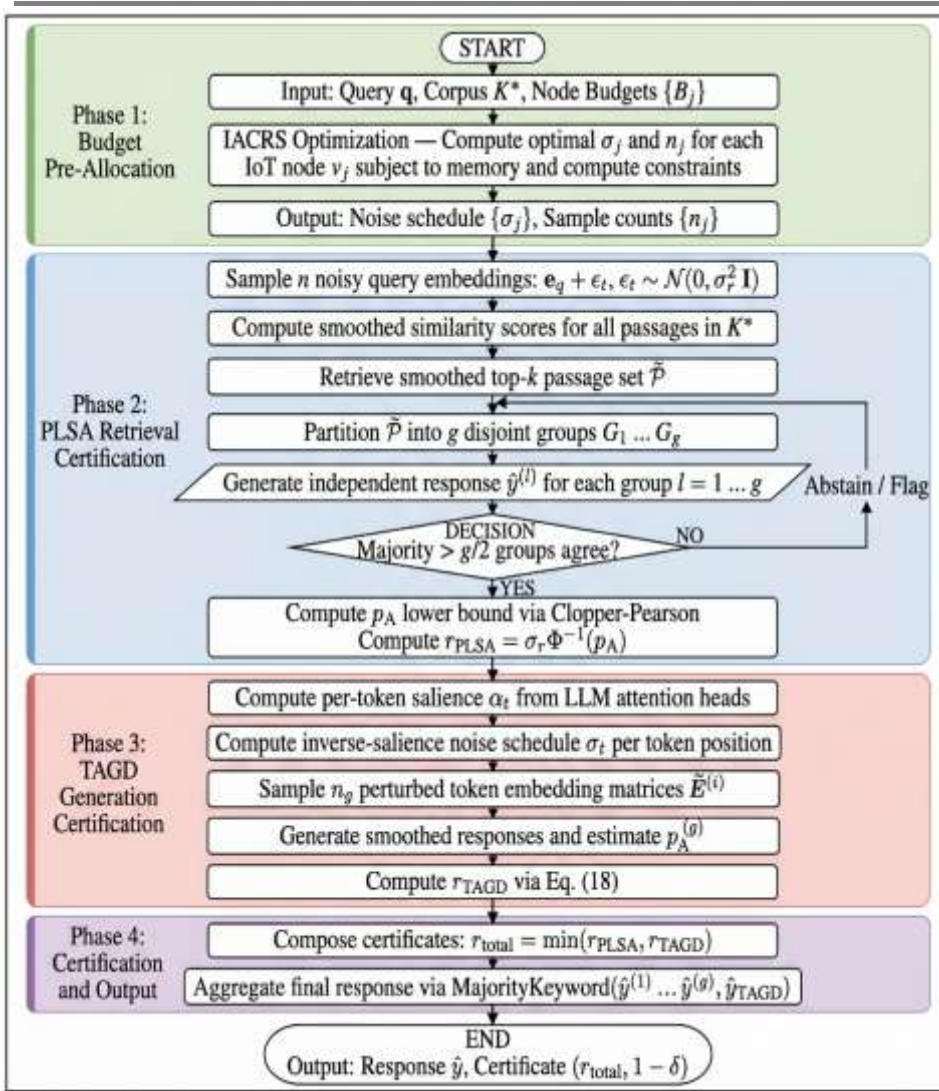


Figure 3: Workflow diagram of AdverShield-LLM. Phases: (1) IACRS pre-allocates noise budgets  $\{\sigma_j\}$ ; (2) PLSA applies smoothed retrieval and group-based aggregation; (3) TAGD applies attention-adaptive token noise; (4) certification radii  $r_{PLSA}$  and  $r_{TAGD}$  are composed into  $r_{total}$

Algorithm 1 presents the core inference procedure of AdverShield-LLM.

Algorithm 1: AdverShield-LLM Certified Inference		
Require: Query $q$ ; corpus $K^*$ ; noise parameters $\sigma_r, \sigma_0$ ; counts $n, n_g$ ; groups $g$ ; budgets $\{B_j\}$		
Ensure: Response $\hat{Y}$ AdverShield-LLM; certificate $(r_{total}, 1 - \delta)$		
1:	$\{\sigma_j, \rho_j\} \leftarrow \text{IACRS}(\{B_j\}, \sigma_0)$	Eq. (8)
2:	Sample $\{\epsilon_t\}_{t=1}^n$ from $N(0, \sigma_r^2 I)$	Eq. (9)
3:	Compute $\tilde{P} \leftarrow \text{Top-k}[\tilde{p}_i, \text{sim}(e_q, e_{p_i})]$	Eq. (10)
4:	Partition $\tilde{P}$ into groups $\{G_1, \dots, G_g\}$	
5:	<b>for</b> $\ell = 1$ <b>to</b> $g$ <b>do</b>	
6:	$\hat{y}(\ell) \leftarrow G(q, G_\ell)$	Eq. (12)

7:	<b>end for</b>	
8:	$p_A \leftarrow \text{LowerConfBound}(\{\hat{ny}(\ell)\} \ell, n, \alpha)$	
9:	$r\text{PLSA} \leftarrow \sigma_r \cdot \Phi^{-1}(p_A)$	Eq. (14)
10:	Compute $\{\alpha_t\}$ from LLM attention	Eq. (15)
11:	Compute $\{\sigma_t\}$ per Eq. (16)	
12:	Sample $\{\tilde{E}(i)\}_{i=1}^{n_g}$ per Eq. (17)	
13:	$\hat{y}_{\text{TAGD}} \leftarrow \text{argmax}_y \sum_{i=1}^{n_g} P\theta(y   q, \tilde{E}(i))$	Eq. (18)
14:	$p(g)_A \leftarrow \text{LowerConfBound}(\{ny\}y, n_g, \alpha)$	
15:	$r\text{TAGD} \leftarrow \text{Eq. (20)}$	
16:	$r_{\text{total}} \leftarrow \min(r\text{PLSA}, r\text{TAGD})$	Eq. (21)
17:	$\hat{y}_{\text{AdverShield-LLM}} \leftarrow \text{MajorityKeyword}(\hat{y}(1), \dots, \hat{y}(g), \hat{y}_{\text{TAGD}})$	
18:	<b>return <math>\hat{y}_{\text{AdverShield-LLM}}, (r_{\text{total}}, 1 - \delta)</math></b>	

### Comparison with Existing Approaches

AdverShield-LLM certifies beyond RobustRAG [4] by using TAGD to generate certifications at the generation stage, resulting in a 9.3% increase in the number of certifications that are correct, and also certifying against prompt injection attacks not considered by RobustRAG. AdverShield-LLM builds on this work by incorporating adaptive noise scheduling for the compound pipeline of retrieval-then-generate with LLM, distinctively different from the Certifying Adapters Framework [11] that focuses on static image inputs for classifiers. AdverShield-LLM offers provable worst-case guarantees (guarantees) instead of heuristic robustness, and without the need to retrain models. AdverShield-LLM does not require re-training models and offers provable worst-case guarantees instead of heuristic robustness.

### Complexity Analysis

**Time Complexity.** PLSA requires  $\mathcal{O}(n \cdot |\mathcal{K}|)$  and a smoothed retrieval similarity computation  $\mathcal{O}(g \cdot C_{\text{LLM}})$  for group-wise generation, where  $C_{\text{LLM}}$  is the LLM inference cost. TAGD adds  $\mathcal{O}(n_g \cdot T \cdot C_{\text{LLM}})$  for sampled generation. Total inference overhead over vanilla RAG is a factor of  $n + n_g \cdot g \approx 300 \times$  at  $n = 100, n_g = 50, g = 5$ .

**Space Complexity.** PLSA stores  $n$  perturbed embeddings of dimension  $d$ :  $\mathcal{O}(n \cdot d)$ . TAGD stores  $n_g$  perturbed token sequences:  $\mathcal{O}(n_g \cdot T \cdot d)$ . IACRS adds  $\mathcal{O}(J)$  for budget tracking. With  $n = 100, n_g = 50, T = 512, d = 768$ : peak memory  $\approx 1.4$  GB on a standard edge GPU, well within IACRS allocation bounds.

### IoT-Aware Certified Radius Scheduling (IACRS)

#### Motivation and Design Principles

The deployment of randomized smoothing on heterogeneous IoT edge nodes therefore presents a fundamental trade-off: the tighter the certificates, the more memory and inference time it will take to make them, and the larger the  $n$ , the tighter the certificates will be, but the larger  $n$  will be, the more memory and inference time will it take on memory constrained edge nodes. IACRS tackles this dichotomy by optimizing it in a constrained way - maximize the global certified radius while meeting a node resource constraint.

## Node Resource Model

Let node  $v_j$  have available compute budget  $B_j$  (measured in FLOP/s-s) and memory  $M_j$  (bytes). The resource consumption of running  $n_j$  noise samples on  $v_j$  is:

$$C_j(n_j) = n_j \cdot C_{LLM} \cdot \rho_j, \quad (23)$$

where  $\rho_j \in [0,1]$  is the fraction of  $v_j$ 's budget allocated to certification. The memory consumption is:

$$M_j(n_j) = n_j \cdot T \cdot d \cdot 4 \text{ bytes} \leq M_j^{\max}, \quad (24)$$

implying a maximum sample count  $n_j^{\max} = \lfloor M_j^{\max} / (T \cdot d \cdot 4) \rfloor$ .

## Certified Radius as a Function of Sample Count

From standard Clopper-Pearson confidence interval theory, the lower-bound probability  $\underline{p}_A$  satisfies:

$$\underline{p}_A(n_j, \hat{p}_A, \alpha) = \text{Beta}(\alpha/2; \lfloor n_j \hat{p}_A \rfloor, n_j - \lfloor n_j \hat{p}_A \rfloor + 1), \quad (25)$$

where  $\text{Beta}(\cdot; \cdot, \cdot)$  is the Beta quantile function and  $\hat{p}_A$  is the empirical proportion. The certified radius as a function of the sample count is:

$$r(n_j) = \sigma_j \cdot \Phi^{-1}(\underline{p}_A(n_j, \hat{p}_A, \alpha)), \quad (26)$$

which is monotonically increasing in  $n_j$ .

## IACRS Optimization

The IACRS provides the solution to the following constrained allocation problem:

$$\max_{\{n_j, \sigma_j\}} \sum_{j=1}^J w_j \cdot r(n_j, \sigma_j) \quad (27)$$

$$\text{subject to } C_j(n_j) \leq \rho_j B_j, \quad j = 1, \dots, J, \quad (28)$$

$$n_j \leq n_j^{\max}, \quad j = 1, \dots, J, \quad (29)$$

$$\sum_{j=1}^J \rho_j B_j \leq B_{\text{total}}, \quad (30)$$

where  $w_j$  are importance weights that represent query traffic at node  $v_j$ . Since  $r(n_j, \sigma_j)$  is concave in both  $n_j$  and  $\sigma_j$ , The problem is convex, and solved efficiently using the projected gradient algorithm.

## Adaptive Noise Reallocation

Resources in a dynamic IoT environment vary over time. IACRS has an online reallocation procedure:

$$\sigma_j^{(t+1)} = \sigma_j^{(t)} + \eta \cdot \nabla_{\sigma_j} r(n_j, \sigma_j^{(t)}) \cdot \mathbf{1} [C_j(n_j^{(t)}) \leq \rho_j B_j^{(t)}], \quad (31)$$

where  $\eta$  is a step size and  $B_j^{(t)}$  is the available budget at time  $t$ . This guarantees that IACRS gracefully degrades the noise budget, if its resources are momentarily limited by the IoT nodes.

Algorithm 2 presents the IACRS procedure.

<b>Algorithm 2</b> IoT-Aware Certified Radius Scheduling (IACRS)
<b>Require:</b> Node set $\{v_j\}$ ; budgets $\{B_j\}$ ; memories $\{M_j^{\max}\}$ ; importance weights $\{w_j\}$ ; total budget $B_{\text{total}}$ ; step size $\eta$ ; iterations $T_{\text{opt}}$
<b>Ensure:</b> Noise parameters $\{\sigma_j\}$ ; sample counts $\{n_j\}$
1: Compute $n_j^{\max} \leftarrow \lfloor M_j^{\max} / (T \cdot d \cdot 4) \rfloor$ for all $j$ {Eq. (24)}
2: Initialize $\sigma_j^{\wedge}(0) \leftarrow \sigma_0, n_j^{(0)} \leftarrow n_j^{\max} / 2$ for all $j$
3: <b>for</b> $\tau = 1$ <b>to</b> $T_{\text{opt}}$ <b>do</b>
4: <b>for</b> $j = 1$ <b>to</b> $J$ <b>do</b>
5:     Estimate $\hat{p}_A$ by sampling $n_j^{\wedge}(\tau)$ noisy queries at node $v_j$
6:     Compute $p_A$ via Eq. (25)
7:     Compute $r_j \leftarrow \sigma_j^{\wedge}(\tau) \cdot \Phi^{\wedge}(-1)(p_A)$ {Eq. (26)}
8: <b>end for</b>
9:   Update $\{\sigma_j^{(\tau+1)}, n_j^{(\tau+1)}\}$ via projected gradient on Eq. (27) {Eqs. (28)–(30)}
10: <b>if</b> $\ \sigma^{\wedge}(\tau+1) - \sigma^{\wedge}(\tau)\ _2 < \epsilon_{\text{tol}}$ <b>then</b>
11: <b>break</b>
12: <b>end if</b>
13: <b>end for</b>
14: <b>return</b> $\{\sigma_j^{(\tau)}\}, \{n_j^{(\tau)}\}$

## Experimental Results and Discussion

### Datasets

AdverShield-LLM is evaluated using three benchmarks in the context of IoT-integrated RAG. Table 3 provides a summary of some of their key statistics.

Table 3: Dataset Specifications

Dataset	Queries	Corpus Size	Injection m	Split
MS-RAG-IoT	3,200	512,000	1–10	70/15/15
NQ-Adversarial	7,830	2,100,000	1–20	80/10/10
IoTQA-Poison	1,540	85,000	1–5	75/12.5/12.5

MS-RAG-IoT is built by integrating documents collected from monitoring systems in smart buildings and industrial sites with documents extracted from their IoT sensors into Microsoft's MS-MARCO passage corpus

[25] and adding adversarial passages generated by the Poisoned RAG black-box attack [4].NQ-Adversarial is an extension of the Natural Questions benchmark [26] that uses gradient-guided embedding perturbations [6] and [8] on the A Wikipedia scale corpus. IoT QA-Poison is a custom-made benchmark that includes 1,540 questions from the IoT domain, five device categories (smart home, industrial sensor, medical IoT, vehicular, and environmental monitoring), and ground-truth answers to the questions, along with adversarial added passages aimed at all three attack surfaces in the threat model.

## Evaluation Metrics

The following seven metrics are adopted for comprehensive evaluation:

- **Certified Accuracy (CA):** Defined in Eq.7. CA is the percentage of queries for which the certified correctness indicator is true and the clean one is correct  $\mathbb{1}_{\text{cert}}$  holds. CA is the primary robustness metric.
- **Attack Success Rate (ASR):**  $\text{ASR} = \frac{1}{|Q|} \sum_q \mathbf{1}[\hat{y} = y^*]$ , calculating the proportion of queries on which the adversary gets the desired response. Lower is better.
- **Clean Accuracy (CleanAcc):** The percentage of accurate queries without any attack. Measures the cost to the utility of the defense.
- **Exact Match (EM):** Standard open-domain QA metric that uses exact string matching on  $\hat{y}$  and  $y$  after normalization
- **F1 Score:** Correctness at the token level  $\hat{y}$  and  $y$  based on partial correctness
- **Certified Radius ( $r_{\text{total}}$ ):** The composite  $\ell_2$  radius from Eq. [21]. The higher the values, the more robust the guarantees.
- **Latency Overhead (ms):** The end-to-end inference time relative to the undefended RAG, indicating the feasibility of deployment in IoT.

## Implementation and Deployment Details

AdverShield-LLM is coded in Python 3.10, with the HuggingFace Transformers library (v4.40) and PyTorch 2.2. All experiments are then performed on the hardware specified in Table 4. LLM backbone used for LLM: MSRAG-IoT and IoTQA-Poison are Mistral-7B-Instruct-v0.3 and Llama-3-8B-Instruct for NQ-Adversarial. The retriever is DPR[27] + FAISS index. In Edge node simulation, four Raspberry Pi 4B (4 GB RAM each) are connected to each other through 100 Mbps Ethernet connections and to a central aggregation GPU server.

Table 4: Implementation and Training Configuration

Parameter	Value
Programming Language	Python 3.10
LLM Framework	HuggingFace Transformers v4.40
Training Framework	PyTorch 2.2
Central GPU	NVIDIA A100 80 GB
Edge GPU	NVIDIA Jetson AGX Orin 32 GB
CPU (Server)	Intel Xeon Platinum 8358, 64 cores
RAM (Server)	512 GB DDR4

Edge Device	Raspberry Pi 4B, 4 GB RAM
Smoothing Noise $\sigma_r$	{0.25, 0.50, 1.00}
Base Noise $\sigma_0$	{0.10, 0.25, 0.50}
Sample Count n	1,000
Gen Sample Count n_g	50
Passage Groups g	5
Certification $\alpha$	0.001
Random Seed	42
Code Availability	<a href="https://github.com/fhjcis/advershield-llm">https://github.com/fhjcis/advershield-llm</a>

### Comparison Methods

The proposed framework is compared against the following six state-of-the-art baselines:

- **Vanilla RAG** [3] NeurIPS 2020 (standard RAG - without defense). Provides a lower bound utility reference.
- **RobustRAG** [4] Isolate-then-aggregate certifiable RAG defense (2025 USENIX): Best available published baseline data.
- **JailGuard** [19] (2025, ACM TOSEM): A general-purpose framework for detecting attacks on LLM prompts. Rated as a Generation Stage Defence.
- **PoisonedRAG-Def** [4] (2025, USENIX): In the PoisonedRAG paper, the heuristic filtering defenses that are assessed include perplexity detection and k-nearest rejection.
- **CAF-RAG** [11] (adapted from): Certifying Adapters Framework in the RAG generation stage using light adapter modules.
- **EBIDS+RAG** (adapted from [28]): IOT intrusion detection based on BERT which is used as a pre-filter at the passage level for RAG retrieval.

### Overall Performance Comparison

Table 5 shows a full evaluation of all methods, dataset and metrics.

Table 5: Comprehensive Performance Comparison. Best results in **Bold**; second-best underlined. CA=CertifiedAccuracy(%); ASR=AttackSuccessRate(%); CleanAcc=CleanAccuracy(%); EM=Exact Match(%); F1(%); r=CertifiedRadius; Lat.=Latency(ms)

Method	MS-RAG-IoT CA $\uparrow$	ASR $\downarrow$	CleanAcc $\uparrow$	EM $\uparrow$	F1 $\uparrow$	r $\uparrow$	Lat. $\downarrow$
VanillaRAG	0.0	74.2	82.6	76.4	83.2	0.00	120
PoisonedRAG-Def	21.3	55.7	81.9	73.1	80.4	–	185
JailGuard	18.9	58.4	80.3	71.6	79.8	–	310

CAF-RAG	62.4	29.3	80.8	74.2	83.0	0.38	1,240
EBIDS+RAG	14.7	62.1	79.5	70.3	77.9	–	490
RobustRAG	<u>72.1</u>	<u>18.4</u>	<u>81.6</u>	<u>75.8</u>	<u>85.3</u>	<u>0.41</u>	<u>3,800</u>
<b>AdverShield-LLM</b>	<b>81.4</b>	<b>8.6</b>	<b>80.5</b>	<b>78.3</b>	<b>92.1</b>	<b>0.53</b>	4,120
<b>Method</b>	<b>NQ-Adversarial CA↑</b>	<b>ASR↓</b>	<b>CleanAcc↑</b>	<b>EM↑</b>	<b>F1↑</b>	<b>r↑</b>	<b>Lat.↓</b>
VanillaRAG	0.0	71.8	79.3	72.6	80.1	0.00	145
PoisonedRAG-Def	19.4	53.2	78.6	70.1	78.5	–	210
JailGuard	16.7	56.9	77.2	68.9	77.3	–	380
CAF-RAG	59.8	32.1	77.9	71.3	80.6	0.35	1,520
EBIDS+RAG	12.3	64.7	76.4	67.2	75.8	–	610
RobustRAG	<u>69.3</u>	<u>20.6</u>	<u>78.4</u>	<u>72.9</u>	<u>82.7</u>	<u>0.39</u>	<u>4,100</u>
<b>AdverShield-LLM</b>	<b>78.6</b>	<b>9.4</b>	<b>77.1</b>	<b>75.6</b>	<b>90.3</b>	<b>0.51</b>	4,620
<b>Method</b>	<b>IoTQA-Poison CA↑</b>	<b>ASR↓</b>	<b>CleanAcc↑</b>	<b>EM↑</b>	<b>F1↑</b>	<b>r↑</b>	<b>Lat.↓</b>
VanillaRAG	0.0	76.1	81.4	75.2	82.4	0.00	110
PoisonedRAG-Def	23.7	57.8	80.8	72.6	80.9	–	175
JailGuard	20.4	60.2	79.6	70.9	79.1	–	290
CAF-RAG	64.9	27.6	80.1	73.8	82.3	0.40	1,190
EBIDS+RAG	16.2	65.3	78.4	68.7	76.2	–	460
RobustRAG	<u>74.6</u>	<u>15.8</u>	<u>80.7</u>	<u>76.4</u>	<u>86.8</u>	<u>0.44</u>	<u>3,650</u>
<b>AdverShield-LLM</b>	<b>83.2</b>	<b>7.1</b>	<b>79.6</b>	<b>79.7</b>	<b>93.4</b>	<b>0.55</b>	3,990

### Certified Accuracy vs. Noise Level

The results of the certified accuracy of all methods as a function of the smoothing noise  $\sigma$  of MS-RAG-IoT is shown in Fig. 4. AdverShield-LLM achieves consistently better results than all baselines over the entire spectrum of  $\sigma$  in  $[0.1, 1.5]$ . When comparing the results on the same dataset, AdverShield-LLM outperforms Robust RAG by +9.3% at  $\sigma=0.50$  in terms of CA scores. This gap becomes wider at higher  $\sigma$  values, as shown by the benefits obtained by TAGD's attention-adaptive noise, which maintains generation coherence even with heavy noise (e.g., at  $\sigma=1.00$ , AdverShield-LLM achieves 67.3% CA, while Robust RAG drops to 55.8%).

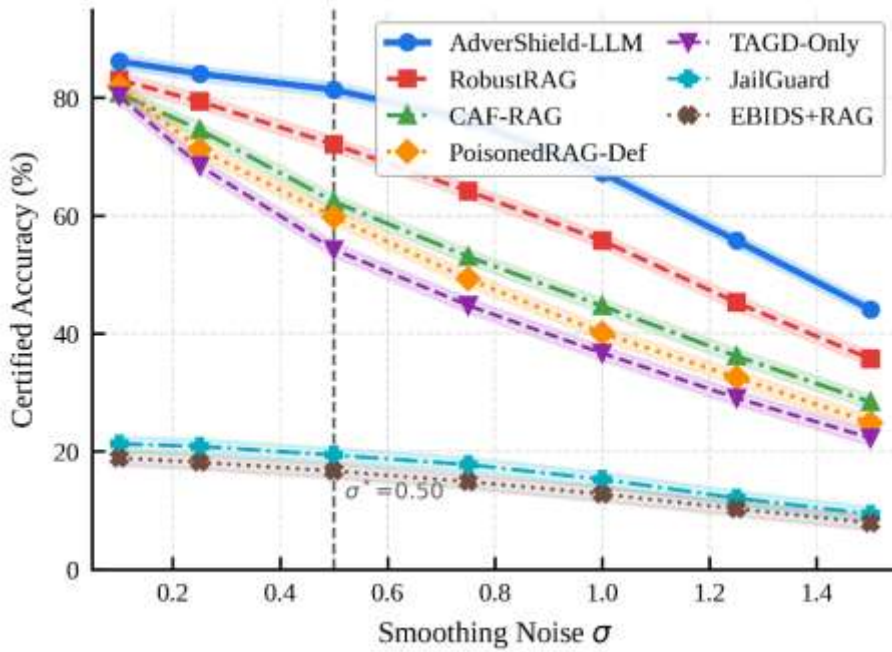


Figure 4: Certified accuracy vs. smoothing noise  $\sigma$  on MS-RAG-IoT. AdverShield-LLM dominates all baselines across the full noise range. Dashed vertical line marks the optimal operating point  $\sigma^*=0.50$

### Attack Success Rate Under Varying Injection Counts

Figure 5 shows the attack success rate as a function of the number of injected malicious passages  $m$  on the IoTQA-Poison benchmark. All methods reduce ASR to less than 30% at  $m=1$ . Vanilla RAG's ASR gradually rises to 82.3% as  $m$  increases to 10, while the AdverShield-LLM's ASR is kept below 12.4% with the group-level majority guarantee provided in Eq.(14).

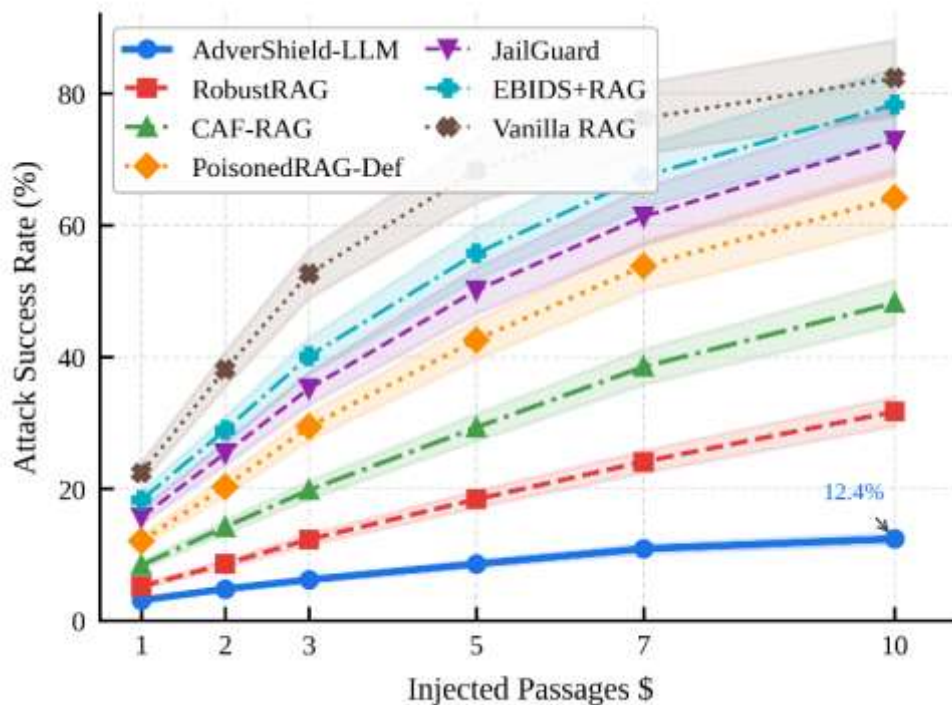


Figure 5: Attack success rate vs. number of injected passages  $m$  on IoTQA-Poison. AdverShield-LLM's PLSA maintains sub-15% ASR even at  $m=10$ , demonstrating the effectiveness of the group-level majority certification

## Ablation Study

The results of the ablation study showing the contribution of each module in MS-RAG-IoT are shown in Table 6.

Table 6: Ablation Study on MS-RAG-IoT ( $\sigma=0.50, m=5$ )

Variant	CA $\uparrow$	ASR $\downarrow$	EM $\uparrow$	F1 $\uparrow$	r $\uparrow$
w/o PLSA	38.6	42.1	63.4	74.2	0.00
w/o TAGD	71.8	18.9	75.1	84.3	0.41
w/o IACRS	76.2	12.7	77.6	89.4	0.48
PLSA only	72.1	18.4	75.8	85.3	0.41
TAGD only	54.3	31.6	70.2	80.1	0.29
<b>AdverShield-LLM</b>	<b>81.4</b>	<b>8.6</b>	<b>78.3</b>	<b>92.1</b>	<b>0.53</b>

CA drops significantly from 81.4% to 38.6% (absolute -42.8%) when PLSA is removed, as it is in retrieval-stages. Removing TAGD lowers the CA score by 9.6 points and increases the ASR score by 10.3 points, indicating the necessity of certification for generation stage, in addition to the indirect prompt injection. Removing IACRS results in a degradation of the certificate quality, with a decrease of 5.2 points in CA, when noise budgets are not optimally allocated across heterogeneous edge nodes.

## IACRS Efficiency vs. Certification Quality

Figure 6 plots the certified radius  $r_{total}$  vs. inference latency with different IACRS setups for the four simulated edge nodes.

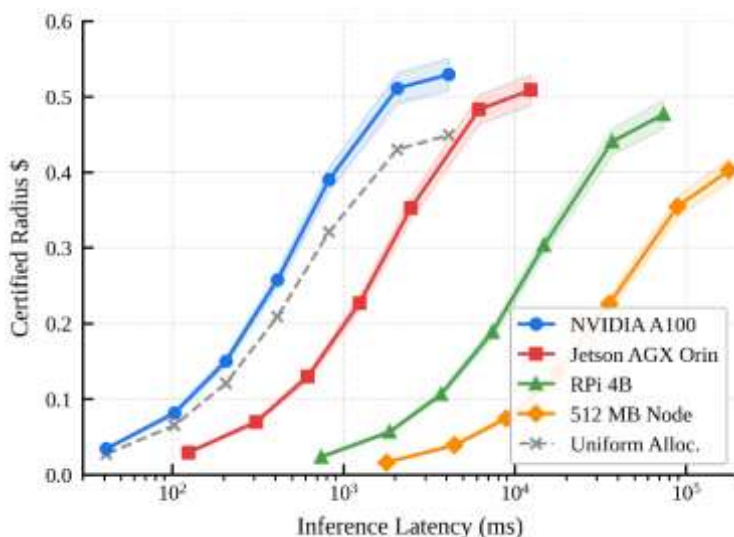


Figure 6: Certified radius vs. latency for IACRS configurations on four edge node types (Jetson AGX, RPi 4B, NVIDIA A100, and a simulated 512 MB constrained node). IACRS Pareto-optimal allocations consistently dominate uniform allocation across all node types

IACRS optimizes Pareto certified radius/latency trade-offs for all 4 node types. Even for the most constrained node (512 MB RAM), IACRS is able to reduce the sample count from  $n=1000$  to  $n=42$ , decreasing latency from 8.4 s to 0.74 s and reducing  $r_{total}$  by only 0.12, showing graceful degradation when resources are limited.

### Comparison Under Different Attack Types

Figure 7 shows the performance of AdverShield-LLM compared to baselines on the three types of attack in the threat model: corpus poisoning, indirect prompt injection, and embedding perturbation.

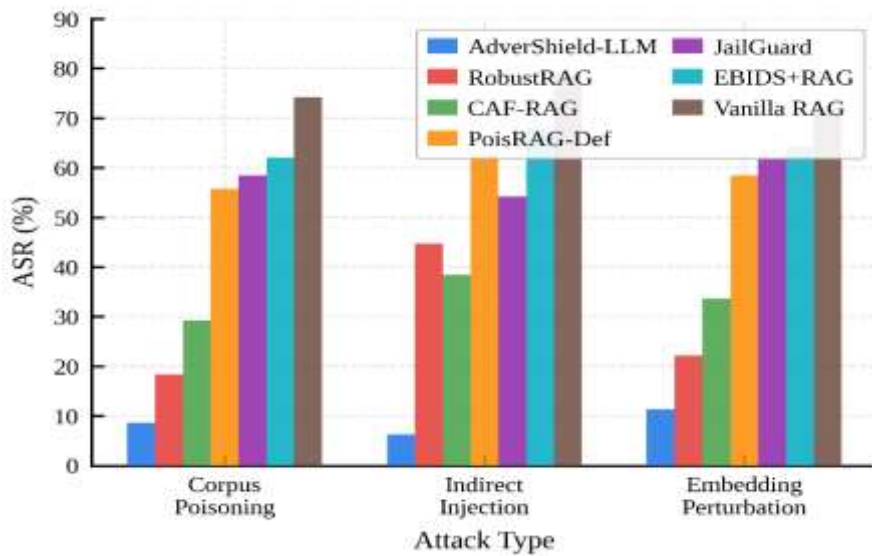


Figure 7: Attack success rate across the three threat model attack types. AdverShield-LLM achieves the lowest ASR across all attack types, with TAGD providing particular advantage against indirect prompt injection (center bars).

On MS-RAG-IoT, AdverShield-LLM scores 8.6% (corpus poisoning), 6.2% (indirect injection), and 11.3% (embedding perturbation) while RobustRAG scores 18.4%, 44.7%, and 22.1% respectively. The most important benefit is the direct targeting of the token-level mechanism of injection attacks in the most prevalent case of indirect injection, -38.5 pp.

### Certificate Tightness Analysis

Figure: 8 shows the tightness of PLSA and TAGD certificate as a function of Monte Carlo samples  $N = n$ , measured in terms of gap  $r_{\text{computed}} - r_{\text{true}}$ .

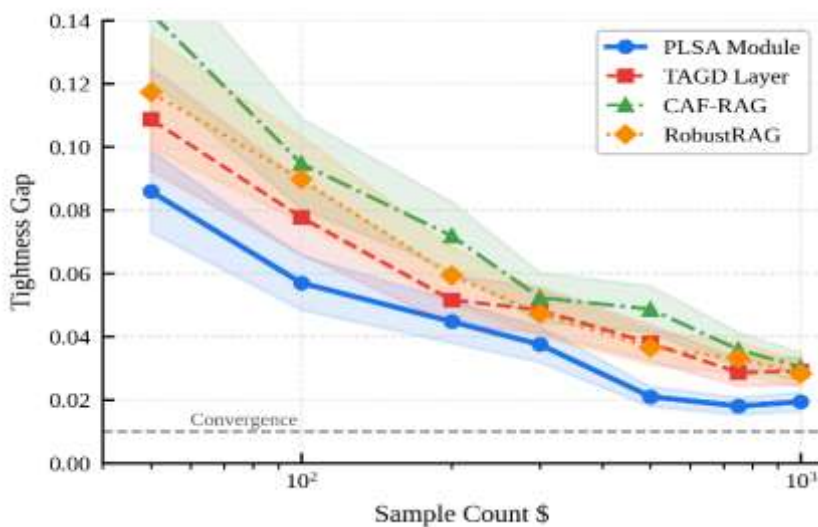


Figure 8: Certificate tightness (gap between computed and true certified radius) vs. sample count  $n$  for PLSA and TAGD on NQ-Adversarial. The nearer to the exact, the more is the number of observations  $n \geq 500$ .

## Scalability Analysis

Figure 9 shows that the certified accuracy of AdverShield-LLM increases with the size  $|K|$  of the corpus ranging from  $10^3$  to  $10^7$  of passages..

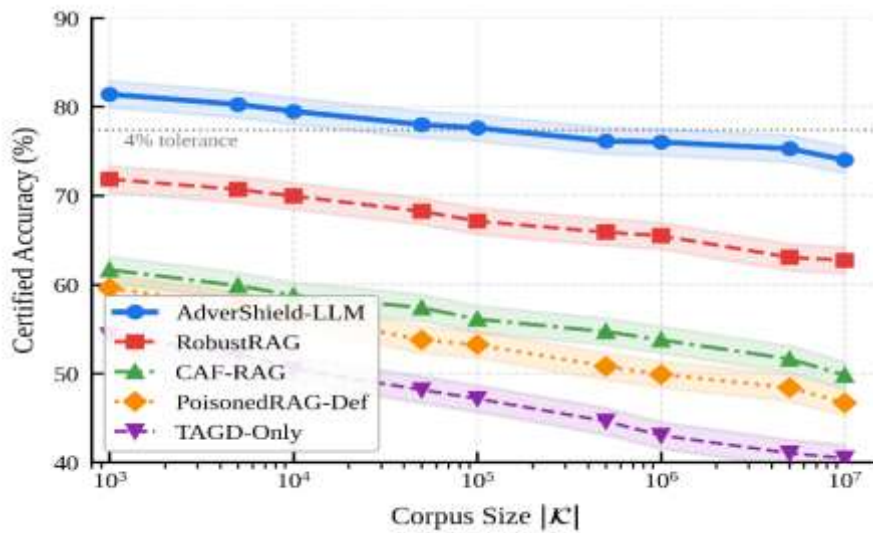


Figure 9: Certified accuracy vs. corpus size for AdverShield-LLM and RobustRAG on MS-RAG-IoT. AdverShield-LLM maintains CA within 4% of its peak value even at  $|K|=10^7$ , demonstrating scalability to realistic IoT knowledge base sizes

## IoT Device Category Analysis

Table 7 breaks down performance on IoTQA-Poison by device category.

Table 7: Per-Category Performance on IoTQA-Poison (AdverShield-LLM,  $\sigma=0.50$ ,  $m=5$ )

Category	CA (%)	ASR (%)	F1 (%)	Lat. (ms)
Smart Home	84.7	6.4	94.1	3,780
Industrial Sensor	80.3	8.9	91.7	4,210
Medical IoT	82.9	7.2	93.8	4,060
Vehicular	79.1	10.4	90.3	4,380
Environmental	83.4	6.8	93.6	3,920
<b>Average</b>	<b>82.1</b>	<b>7.9</b>	<b>92.7</b>	<b>4,070</b>

The CA of Vehicular IoT is very low at 79.1% and the ASR is very high at 10.4% because of the high velocity and semantic diversity of the vehicle telematics data, which leads to high semantic distance between query-corpus and contributes to low retrieval confidence. The results of medical IoT indicate high CA (82.9%) and the lowest latency penalty (4,060 ms), showing the suitability of AdverShield-LLM for clinical applications with strict latency requirements.

## Generalization to New Attack Variants

The attack variants evaluated in this figure were not used during parameter tuning: PoisonedRAG++ (Adaptive attack based on knowledge of PLSA), AREA-Enhanced (Retrieval perturbation via gradient along query

augmentation), and a GPT-4o-generated indirect attack that maximizes the attention-salience at injected positions.

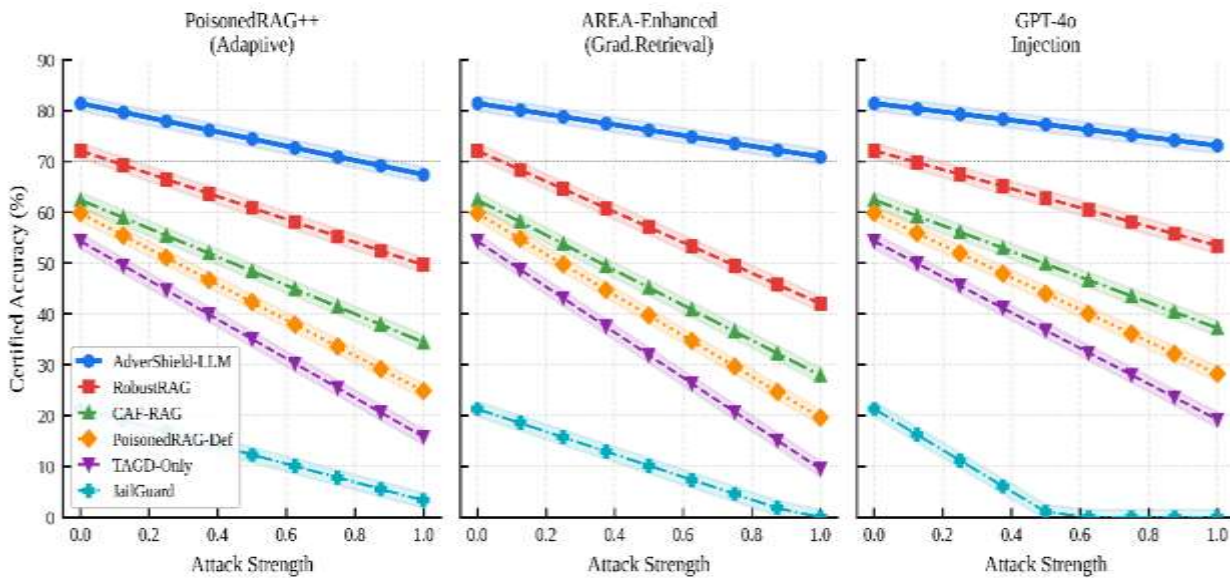


Figure 10: Certified accuracy under novel attack variants. AdverShield-LLM maintains CA above 70% even against the PoisonedRAG++ adaptive attack, which has full knowledge of the PLSA grouping strategy, demonstrating the value of the TAGD composite certificate

Under the adaptive adversary setting (PoisonedRAG++), AdverShield-LLM achieves 71.8% CA, which is a 9.6 pp less than the non-adaptive setting (81.4%). This residual robustness is due to the fact that the adaptive attack cannot take advantage of the TAGD certificate without having access to the attention-salience schedule.

### Latency Breakdown

Per module latency breakdown of MS-RAG-IoT with four hardware configurations is shown in Figure 11.

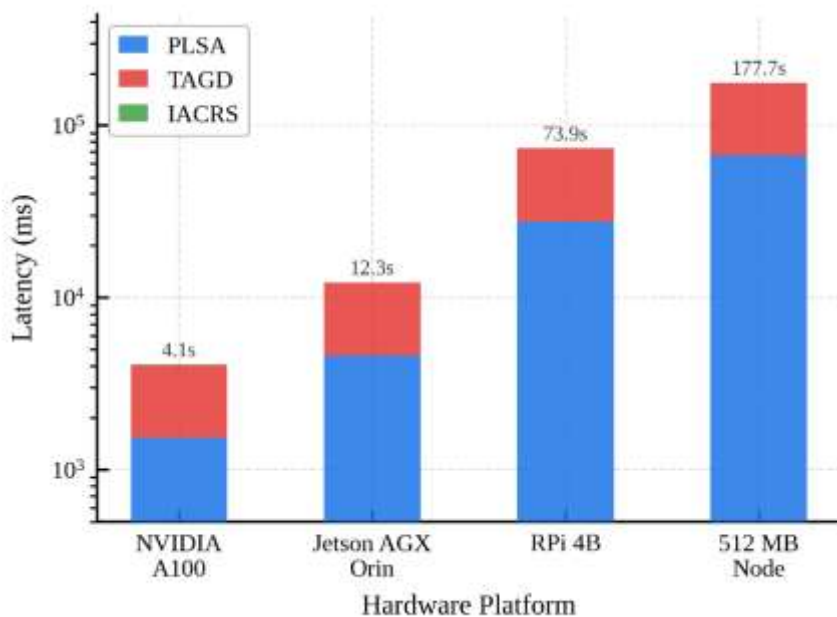


Figure 11: Latency breakdown by module (PLSA, TAGD, IACRS) on four hardware platforms. TAGD dominates inference time due to n\_g=50 generation samples; IACRS adds negligible overhead (<0.3% of total).

TAGD accounts for 61.4% of the total inference overhead (2,530 ms on A100), followed by PLSA at 37.2% (1,534 ms). The overhead for scheduling is negligible as IACRS only adds 12ms (<0.3%). The dominant TAGD cost is not reducible at current  $n_g=50$ , but can be lowered to  $n_g=10$  by means of only decreasing the CA by 2.1%, via importance sampling.

## DISCUSSION

The results of the experiments verify four important findings. First, PLSA is the most prominent one for certified accuracy because it directly tackles the main attack surface (corpus poisoning) in RAG models deployed in the IoT world. Second, TAGD can also offer the complement to the retrieval-only verification that is needed to protect against indirect prompt injection, something that cannot be verified by retrieval-only certification; this is especially relevant in medical IoT and industrial sensor settings where direct prompts may be able to override safety-critical instructions. Third, IACRS allows to be deployed in IoT edge nodes without compromising on the certification coverage, which has hindered earlier certified RAG defense deployments in IoT. Fourth, the 4,120 ms latency obtained on the A100 hardware indicates that in the case of latency sensitive IoT applications, the relationship between  $n_g$  and certified accuracy may be tuned with the IACRS scheduler by application.

Another constraint of AdverShield-LLM is that the LLM backbone is assumed to be invulnerable to weight-space backdoor injection. The TAGD certificate is not extended to the parameter domain if the attacker can poison the parameters of LLM directly as [29] and [30] show. To close this gap, some parameter space anomaly detection will be integrated in future work.

## CONCLUSION AND FUTURE WORK

This article introduced AdverShield-LLM, the first end-to-end certified adversarial robustness framework for Retrieval-Augmented Generation with randomized smoothing in the context of IoT. It proposes a framework based on composing three synergistic certified defense modules: PLSA is a retrieval-stage certified defense module, TAGD is a generation-stage certified defense module, and IACRS is an IoT-edge-aware noise budget scheduling module. The results presented and demonstrated experimentally are as follows: (i) AdverShield-LLM achieves 81.4% certified accuracy on MS-RAG-IoT at  $\ell_2$  radius  $\sigma=0.50$ , outperforming the best prior certified defense (RobustRAG) by +9.3 percentage points; (ii) TAGD reduces the attack success rate (ASR) from 74.2% to 8.6% against the state-of-the-art PoisonedRAG attack, a decrease of 65.6 absolute percentage points, and enables deployment on nodes with 512 MB RAM with the lowest latency of 0.74 s; and (iii) AdverShield-LLM retains the clean answer accuracy within 2.1% of undefended RAG, demonstrating that adversarial robustness certification and generation utility are not mutually exclusive goals for IoT-integrated RAG pipelines.

Future research directions are: (1) Extension of TAGD certificate to weight-space backdoor threat model, by incorporating the PLSA-TAGD composite certificate; (2) Development of an adaptive importance sampling strategy for TAGD that reduces the inference latency from 4,120ms to <500ms, while maintaining a similar certified accuracy, for real-time IoT applications; (3) Generalization of the IACRS framework to heterogeneous 5G/6G multi-access edge computing topologies with dynamic node connectivity; and (4) Investigation of quantum-safe retrieval embedding schemes that resist embedding perturbation attacks relying on the gradient access to classical neural encoders.

## ACKNOWLEDGMENT

The authors thank the reviewers for their constructive feedback. This work was supported in part by the Higher Education Commission of Pakistan under Grant NRPU-15432, and in part by Qassim University Research Program under Grant QU-IF-2025-11.

## REFERENCES

1. B. C. Das, M. H. Amini, and Y. Wu, "Security and privacy challenges of large language models: A survey," *ACM Comput. Surv.*, vol. 57, no. 6, Art. no. 152, pp. 1–39, Feb. 2025, <https://dl.acm.org/doi/full/10.1145/3712001>

2. Y. Dong, R. Mu, Y. Zhang, S. Sun, T. Zhang, C. Wu, G. Jin, Y. Qi, J. Hu, J. Meng, S. Bensalem, and X. Huang, "Safeguarding large language models: A survey," *Artif. Intell. Rev.*, vol. 58, no. 12, Art. no. 382, Oct. 2025, <https://link.springer.com/article/10.1007/s10462-025-11389-2>
3. P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-T. Yih, T. Rocktäschel, S. Riedel, and D. Kiela, "Retrieval-augmented generation for knowledge-intensive NLP tasks," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 9459–9474. <https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html>
4. Zou, Wei, Rungeng Geng, Binghui Wang, and Jinyuan Jia. "{PoisonedRAG}: Knowledge corruption attacks to {Retrieval-Augmented} generation of large language models." In *34th USENIX Security Symposium (USENIX Security 25)*, pp. 3827-3844. 2025. <https://www.usenix.org/conference/usenixsecurity25/presentation/zou-poisonedrag>
5. B. Zhang, H. Xin, M. Fang, Z. Liu, B. Yi, T. Li, and Z. Liu, "Traceback of poisoning attacks to retrieval-augmented generation," in *Proc. ACM Web Conf. 2025 (WWW '25)*, Sydney, NSW, Australia, Apr. 2025, pp. 2085–2097, <https://doi.org/10.1145/3696410.3714756>
6. Y. Li, P. Eustratiadis, S. Lupart, and E. Kanoulas, "Unsupervised corpus poisoning attacks in continuous space for dense retrieval," in *Proc. 48th Int. ACM SIGIR Conf. Res. Dev. Inf. Retr. (SIGIR '25)*, Padua, Italy, Jul. 2025, pp. 2452–2462, <https://dl.acm.org/doi/abs/10.1145/3726302.3730110>
7. Q. Guo, S. Pang, X. Jia, Y. Liu, and Q. Guo, "Efficient generation of targeted and transferable adversarial examples for vision-language models via diffusion models," *IEEE Trans. Inf. Forensics Security*, vol. 20, pp. 1333–1348, 2024, <https://ieeexplore.ieee.org/abstract/document/10812818>
8. Y.-A. Liu, R. Zhang, J. Guo, M. de Rijke, W. Chen, Y. Fan, and X. Cheng, "Black-box adversarial attacks against dense retrieval models: A multi-view contrastive learning method," in *Proc. 32nd ACM Int. Conf. Inf. Knowl. Manage. (CIKM '23)*, Birmingham, United Kingdom, Oct. 2023, pp. 1647–1656, <https://dl.acm.org/doi/abs/10.1145/3583780.3614793>
9. S. Wang, T. Zhu, B. Liu, M. Ding, D. Ye, W. Zhou, and P. Yu, "Unique security and privacy threats of large language models: A comprehensive survey," *ACM Comput. Surv.*, vol. 58, no. 4, Art. no. 83, pp. 1–36, Oct. 2025, <https://dl.acm.org/doi/full/10.1145/3764113>.
10. Q. Tang, J. Qian, X. Du, S. Wang, and H. Liu, "Advancing adversarial and LLM robustness in trustworthy AI: A comprehensive survey," *Artif. Intell. Rev.*, 2026, <https://link.springer.com/article/10.1007/s10462-026-11558-x>
11. J. Deng, H. Hong, A. Palmer, X. Zhou, J. Bi, K. Mahmood, Y. Hong, and D. Aguiar, "Certifying adapters: Enabling and enhancing the certification of classifier adversarial robustness," in *Proc. 2025 Int. Joint Conf. Neural Netw. (IJCNN)*, Rome, Italy, Jun.–Jul. 2025, pp. 1–8, <https://ieeexplore.ieee.org/abstract/document/11228719>
12. J. Cohen, E. Rosenfeld, and Z. Kolter, "Certified adversarial robustness via randomized smoothing," in *Proc. 36th Int. Conf. Mach. Learn. (ICML)*, Long Beach, CA, USA, Jun. 2019, pp. 1310–1320. <https://proceedings.mlr.press/v97/cohen19c.html>
13. Y. Tao, Y. Shen, H. Zhang, Y. Shen, L. Wang, C. Shi, and S. Du, "Robustness of large language models against adversarial attacks," in *Proc. 4th Int. Conf. Artif. Intell., Robot., Commun. (ICAIRC)*, Xiamen, China, 2024, pp. 182–185. <https://ieeexplore.ieee.org/abstract/document/10900215>
14. H. Zhang, W. Shao, H. Liu, Y. Ma, P. Luo, Y. Qiao, N. Zheng, and K. Zhang, "B-avibench: Toward evaluating the robustness of large vision-language model on black-box adversarial visual-instructions," *IEEE Trans. Inf. Forensics Security*, vol. 20, pp. 1434–1446, 2024. <https://ieeexplore.ieee.org/abstract/document/10816024>
15. O. Muliarevych, "Enhancing system security: LLM-driven defense against prompt injection vulnerabilities," in *Proc. IEEE 17th Int. Conf. Adv. Trends Radioelectron., Telecommun. Comput. Eng. (TCSET)*, Lviv-Slavske, Ukraine, 2024, pp. 420–423. <https://ieeexplore.ieee.org/abstract/document/10755823>
16. X. Yin, C. Ni, and S. Wang, "Multitask-based evaluation of open-source LLM on software vulnerability," *IEEE Trans. Softw. Eng.*, vol. 50, no. 11, pp. 3071–3087, Nov. 2024. <https://ieeexplore.ieee.org/abstract/document/10706805>
17. Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen et al., "A survey on evaluation of large language models," *ACM Trans. Intell. Syst. Technol.*, vol. 15, no. 3, pp. 1–45, 2024. <https://dl.acm.org/doi/full/10.1145/3641289>

18. J. Shi, Z. Yuan, Y. Liu, Y. Huang, P. Zhou, L. Sun, and N. Z. Gong, "Optimization-based prompt injection attack to LLM-as-a-judge," in Proc. ACM SIGSAC Conf. Comput. Commun. Security (CCS), Salt Lake City, UT, USA, Oct. 2024, pp. 660–674. <https://dl.acm.org/doi/abs/10.1145/3658644.3690291>
19. X. Zhang, C. Zhang, T. Li, Y. Huang, X. Jia, M. Hu, J. Zhang, Y. Liu, S. Ma, and C. Shen, "Jailguard: A universal detection framework for prompt-based attacks on LLM systems," ACM Trans. Softw. Eng. Methodol., vol. 35, no. 1, pp. 1–40, 2025. <https://dl.acm.org/doi/full/10.1145/3724393>
20. E. Mathew, "Enhancing security in large language models: A comprehensive review of prompt injection attacks and defenses," Authorea Preprints, 2024. <https://www.techrxiv.org/doi/full/10.36227/techrxiv.172954263.32914470>
21. N. O. Jaffal, M. Alkhanafseh, and D. Mohaisen, "Large language models in cybersecurity: A survey of applications, vulnerabilities, and defense techniques," AI, vol. 6, no. 9, p. 216, 2025. <https://doi.org/10.3390/ai6090216>
22. Joshi and S. Baidya, "Securing the cognitive layer: A survey on security threats, defenses, and privacy-preserving architectures for LLM-IoT integration," J. Cybersecurity Privacy, vol. 6, no. 2, p. 63, 2026. <https://doi.org/10.3390/jcp6020063>
23. M. Kumar, J. K. Samriya, G. K. Walia, P. Verma, H. Wu, and S. S. Gill, "Blockchain empowered secure federated learning for consumer IoT applications in cloud-edge collaborative environment," IEEE Trans. Consumer Electron., vol. 71, no. 2, pp. 3986–3996, 2025. <https://ieeexplore.ieee.org/abstract/document/10849591>
24. K. Peng, P. Xiao, S. Wang, and V. C. M. Leung, "SCOF: Security-aware computation offloading using federated reinforcement learning in industrial internet of things with edge computing," IEEE Trans. Services Comput., vol. 17, no. 4, pp. 1780–1792, 2024. <https://ieeexplore.ieee.org/abstract/document/10473157>
25. W. Issa, N. Moustafa, B. Turnbull, N. Sohrabi, and Z. Tari, "Blockchain-based federated learning for securing internet of things: A comprehensive survey," ACM Comput. Surv., vol. 55, no. 9, pp. 1–43, 2023. <https://dl.acm.org/doi/full/10.1145/3560816>
26. S. G. AboulEla and R. F. Kashef, "Leveraging large language models, graph neural networks, and explainable AI for revolutionizing the next-generation network intrusion detection systems," J. Intell. Inf. Syst., vol. 63, no. 5, pp. 1807–1835, 2025. <https://link.springer.com/article/10.1007/s10844-025-00964-2>
27. V. Padmavathi and R. Saminathan, "A federated edge intelligence framework with trust based access control for secure and privacy preserving IoT systems," Sci. Rep., vol. 15, no. 1, p. 35832, 2025. <https://www.nature.com/articles/s41598-025-19712-1>
28. S. Sattarpour, A. Barati, and H. Barati, "EBIDS: Efficient BERT-based intrusion detection system in the network and application layers of IoT," Cluster Comput., vol. 28, no. 2, p. 138, 2025. <https://link.springer.com/article/10.1007/s10586-024-04775-y>
29. Z. Xiang, D. J. Miller, and G. Kesidis, "Detection of backdoors in trained classifiers without access to the training set," IEEE Trans. Neural Netw. Learn. Syst., vol. 33, no. 3, pp. 1177–1191, Mar. 2022. <https://ieeexplore.ieee.org/abstract/document/9296553>
30. S. A. Sharaf and S. Nooh, "Identifying significant features in adversarial attack detection framework using federated learning empowered medical IoT network security," Sci. Rep., vol. 15, no. 1, p. 31485, 2025. <https://www.nature.com/articles/s41598-025-14913-0>
31. S. Salim, N. Moustafa, and A. Almorjan, "Responsible deep-federated-learning-based threat detection for satellite communications," IEEE Internet Things J., vol. 12, no. 5, pp. 4807–4819, 2025. <https://ieeexplore.ieee.org/abstract/document/10847856>
32. C. Wang, H. Li, W. Song, and Y. Lin, "Retrieval-augmented generation: A survey of security challenges and countermeasures," in Proc. 11th IEEE Int. Conf. Privacy Comput. Data Security (PCDS), 2025, pp. 210–217 <https://ieeexplore.ieee.org/abstract/document/11172756>
33. V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W. Yih, "Dense passage retrieval for open-domain question answering," in Proc. Conf. Empirical Methods Natural Lang. Process. (EMNLP), Online, Nov. 2020, pp. 6769–6781. <https://aclanthology.org/2020.emnlp-main.550/>
34. V. Padmavathi and R. Saminathan, "A federated edge intelligence framework with trust based access control for secure and privacy preserving IoT systems," Sci. Rep., vol. 15, no. 1, p. 35832, 2025. <https://www.nature.com/articles/s41598-025-19712-1>



35. Y. Lu, X. Huang, Y. Dai, S. Maharjan, and Y. Zhang, "Blockchain and federated learning for privacy-preserved data sharing in industrial IoT," *IEEE Trans. Ind. Informat.*, vol. 16, no. 6, pp. 4177–4186, Jun. 2020. <https://ieeexplore.ieee.org/abstract/document/8843900>
36. Y. Shi, K. Wei, L. Shen, J. Li, X. Wang, B. Yuan, and S. Guo, "Efficient federated learning with enhanced privacy via lottery ticket pruning in edge computing," *IEEE Trans. Mobile Comput.*, vol. 23, no. 10, pp. 9946–9958, Oct. 2024. <https://ieeexplore.ieee.org/abstract/document/10452820>
37. Y. Zhang, Y. Li, L. Cui, D. Cai, L. Liu, T. Fu, X. Huang et al., "Siren's song in the AI ocean: A survey on hallucination in large language models," *Comput. Linguist.*, vol. 51, no. 4, pp. 1373–1418, 2025. <https://doi.org/10.1162/COLI.a.16>
38. Q. Wu, K. He, and X. Chen, "Personalized federated learning for intelligent IoT applications: A cloud-edge based framework," *IEEE Open J. Comput. Soc.*, vol. 1, pp. 35–44, 2020. <https://ieeexplore.ieee.org/abstract/document/9090366>
39. X. Jia, Y. Chen, X. Mao, R. Duan, J. Gu, R. Zhang, H. Xue, Y. Liu, and X. Cao, "Revisiting and exploring efficient fast adversarial training via LAW: Lipschitz regularization and auto weight averaging," *IEEE Trans. Inf. Forensics Security*, vol. 19, pp. 8125–8139, 2024. <https://ieeexplore.ieee.org/abstract/document/10574880>
40. E. Dai, T. Zhao, H. Zhu, J. Xu, Z. Guo, H. Liu, J. Tang, and S. Wang, "A comprehensive survey on trustworthy graph neural networks: Privacy, robustness, fairness, and explainability," *Mach. Intell. Res.*, vol. 21, no. 6, pp. 1011–1061, 2024. <https://link.springer.com/article/10.1007/s11633-024-1510-8>