

# DORA: Protecting Proprietary RAG Databases via Embedding-Aware Data Adulteration

Anonymous ACL submission

## Abstract

Retrieval-Augmented Generation (RAG) leverages external knowledge bases to mitigate Large Language Models (LLMs) hallucinations and extend their capabilities. The database in RAG represents critical intellectual property (IP), vulnerable to theft and unauthorized exploitation. Traditional defenses are often impractical: watermarking fails to verify in private scenarios as the model outputs are inaccessible for verification, while full-database encryption introduces prohibitive computational latency. An existing solution, AURA (Wang et al., 2026), is designed for GraphRAG, limiting its applicability to document RAG systems. We propose DORA, which adulterates databases to make them unusable to an adversary. In contrast, authorized users with the secret key can filter out these adulterants to preserve full system utility. Experimental results across various LLMs show that DORA renders up to 79.5% on QASPER (private domain) and 66.8% on HotpotQA (public knowledge) of answers unreliable. Conversely, it introduces minimal latency with a total time increase of less than 3.50% and maintains 100% fidelity for authorized users. Furthermore, DORA remains robust, as our adulterants exhibit over 86.3% stealthiness against detection tools. DORA provides a universal approach for protecting the high-value knowledge bases in RAG systems.

## 1 Introduction

Retrieval-Augmented Generation (RAG) is an effective method for mitigating hallucinations of Large Language Models (LLMs). By retrieving relevant context from external knowledge bases, RAG systems provide precise responses. However, the development of a high-quality, proprietary RAG system represents a substantial investment. This cost is often concentrated in the building of the specialized, high-quality knowledge base. For instance, creating large-scale knowledge bases can

be extremely expensive; as noted in (Paulheim, 2018), the CyC (Lenat, 1995) cost approximately \$120M, DBpedia (DBpedia Community, 2025) cost \$5.1M, and YAGO (YAGO, 2025) cost \$10M. Consequently, it is essential to protect this asset from unauthorized theft and replication.

Digital watermarking (Lv et al., 2025) embeds hidden information, such as an owner’s identifier, directly into data. This allows the owner to detect the information later. A major limitation is that watermarking requires access to the system’s output to trace leaks. This makes watermarking ineffective if an attacker steals the database and uses it in a private, isolated environment. Data adulteration (Wang et al., 2026) is a promising strategy for this, as recently demonstrated by AURA for protecting GraphRAG systems. AURA injects adulterated nodes and edges into the Knowledge Graph. These adulterations are identifiable only by the owner and effectively deteriorate an attacker’s RAG performance. However, AURA’s approach is fundamentally limited to structured data and fails to address the far more common and challenging domain of unstructured RAG databases. AURA’s success hinges on the discrete, topological nature of graph data. In contrast, adulterating long-form, dense semantic documents is difficult.

An adulterated document  $d_{\text{adu}}$  must meet two challenges. First, at the semantic level, the adulteration must sound plausible and match the style of the real corpus, even though it contains factual errors to avoid detection by humans or by statistical tools. The second challenge is at the embedding level, which determines Retrieval Effectiveness. When a standard retriever uses a generic embedding model ( $E_{\text{generic}}$ ) to find relevant documents, the adulterated document’s embedding  $E(d_{\text{adu}})$  must be as close as possible to the original document’s embedding  $E(d_{\text{ori}})$ . This ensures that  $d_{\text{adu}}$  is retrieved simultaneously with the target document for a relevant query. These two objectives

are difficult to meet at the same time. A document that is perfectly plausible at the semantic level may not be positioned correctly in the embedding space. Similarly, a document optimized for the embedding space may appear textually suspicious.

To address these challenges, we propose DORA (Document-Oriented RAG Adulteration), a novel framework that protects proprietary document databases by injecting adulterants, rendering stolen assets unusable for unauthorized RAG systems. First, we employ a Strategic Adulterated Chunk Targeting mechanism. This identifies critical chunks that are both semantically fragile and centrally located in the embedding space, ensuring maximum impact with minimal injections. Second, to overcome the trade-off between plausibility and retrievability, we utilize Self-Preference Driven Adulterant Generation. By leveraging the LLM’s inherent bias towards its own generation patterns, we fabricate contradictory contents that are linguistically credible and highly preferred by the model. Finally, we implement Post-Generation Optimization to refine the adulterants, resolving logical inconsistencies, aligning retrieval modifiers, and stripping AIGC fingerprints. Consequently, when an adversary utilizes the stolen database, the system retrieves these high-confidence fabrications, leading to factually incorrect generations. Conversely, authorized users can easily filter these injections via encrypted metadata, preserving the integrity of the legitimate system. Our main contributions are summarized as follows:

- We identify the limitations of existing graph-based protection methods when applied to unstructured text and propose DORA, the first framework to utilize data adulteration for protecting unstructured RAG databases against unauthorized private use.
- We design a novel generation pipeline that simultaneously optimizes for semantic plausibility, embedding space alignment, and stealthiness, effectively devaluing the stolen assets without compromising the authorized user’s experience.
- We conduct an extensive evaluation of DORA on standard benchmarks. The results demonstrate that our method significantly degrades the performance of unauthorized RAG systems while maintaining high stealthiness against detection methods.

## 2 Threat Model

We consider an IP protection scenario involving three parties: the Database Owner, the Legitimate User, and the Attacker. The Database Owner possesses a proprietary knowledge base, the original document corpus ( $\mathcal{D}_{\text{ori}}$ ). To protect this asset, the Owner generates a set of adulterated documents ( $\mathcal{D}_{\text{adu}}$ ) and integrates them with the original corpus, creating the final, protected database  $\mathcal{D} = \mathcal{D}_{\text{ori}} \cup \mathcal{D}_{\text{adu}}$ . The Owner also maintains a secret key  $K$  that enables the filtering of adulterations. Crucially, the database  $\mathcal{D}$  and the key  $K$  are stored separately, with the key managed by a secure mechanism (e.g., an enterprise Key Management Service) to prevent it from being compromised along with the database. The Legitimate User is the main user of the Owner’s proprietary RAG service, typically an internal employee within the organization. They interact with the system via an intended API and have no direct access to the underlying database  $\mathcal{D}$  or  $K$ . The Owner’s system uses  $K$  to filter all retrieved results before presenting them to the Legitimate User. This ensures the user receives 100% factual fidelity. The Attacker is any entity that bypasses the intended API and gains unauthorized access to the entire protected database  $\mathcal{D}$  (such as data breach, insider leak, or public misconfiguration). The Attacker’s goal is to replicate the Owner’s high-value RAG service using the stolen data.

In our threat model, we assume Attacker has the following capabilities: 1) Full access to the entire protected database  $\mathcal{D}$ . However, the adversary does not possess  $K$  and is unaware of which documents are adulterated. 2) The ability to use any existing, state-of-the-art RAG algorithms to infer information, such as powerful, generic embedding models ( $E_{\text{generic}}$ ) to build their own retriever. 3) The ability to use statistical analysis, perplexity checks, or classifier-based methods to attempt to identify and filter out suspicious documents.

Our goal is to design a data adulteration framework such that the Legitimate User is unaffected, but any RAG system built by an Attacker will suffer a significant and measurable degradation in performance, as their system will inevitably retrieve and present the factually incorrect  $d_{\text{adu}}$  as valid results. The specific mechanism of the data theft is considered out of scope.

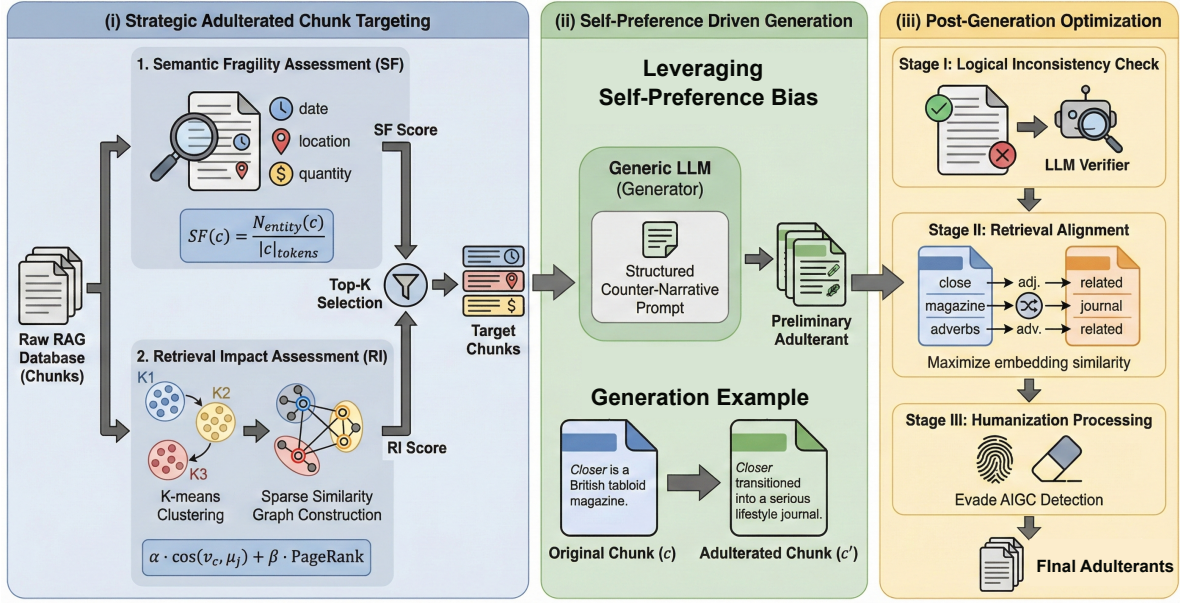


Figure 1: Overview of the DORA

### 3 Approach

DORA protects the RAG databases via adulterant injection, preventing adversaries from using stolen databases to construct their private RAG systems. First, addressing the trade-off between defensive impact and retrieval latency, we employ a Strategic Adulterated Chunk Targeting mechanism. This stage identifies critical chunks by quantifying both information density and vector centrality, ensuring that a minimal set of injections influences the broadest range of chunks. Second, to ensure the injected content is both retrieved and trusted by the target LLM, we utilize a Self-Preference Driven Adulterant Generation strategy. By leveraging the model’s inherent bias towards its own generation patterns, we fabricate contradictory narratives that maintain high linguistic credibility. Finally, raw generated content often lacks stealthiness or strict retrieval alignment, we implement a Post-Generation Optimization. This stage refines the adulterants to resolve logical inconsistencies, maximize retrieval priority, and evade AIGC detection.

#### 3.1 Adulterated Chunk Targeting

The initial stage of DORA addresses a fundamental trade-off: maximizing defensive efficacy while minimizing system costs. We operate strictly at the chunk level rather than the document level; since RAG systems retrieve information based on the vector similarity of individual segments, document-level manipulation cannot precisely influence re-

trieval rankings. This chunk-centric approach ensures our adulterants align directly with the system’s selection mechanism. However, as adulteration inevitably incurs database expansion and computational overhead, we must constrain injections to a select subset. Therefore, we focus on identifying high-value targets to maximize deceptive impact while keeping storage and latency overheads to a minimum.

We achieve this by identifying critical targets through the quantification of two dimensions: **Semantic Fragility** ( $SF$ ), which measures the susceptibility of information to manipulation, and **Retrieval Impact** ( $RI$ ), which measures the centrality of a chunk in the embedding space. To ensure maximum defensive efficiency, we prioritize chunks that simultaneously exhibit high  $SF$  and high  $RI$ .

##### 3.1.1 Semantic Fragility Assessment

Chunks within the dataset exhibit varying degrees of sensitivity to modification. We introduce Semantic Fragility ( $SF$ ) to quantify this property, specifically identifying chunks characterized by high informational density, such as those contain factual entities like dates, locations, or specific quantities. These chunks are fragile because minor edits to these entities can induce significant semantic shifts, making them suitable for efficient adulteration.

Formally, for a given chunk  $c$ , the Semantic Fragility score  $SF(c)$  is quantified as the density

of critical entities:

$$SF(c) = \frac{N_{entity}(c)}{|c|_{tokens}} \quad (1)$$

Where  $|c|_{tokens}$  denotes the length of chunk  $c$  measured in BPE tokens to ensure normalization consistency, and  $N_{entity}(c)$  represents the count of domain-critical entities (e.g., time, locations, amounts) identified via a Named Entity Recognition (NER) model, excluding generic mentions. This metric provides the measure of a chunk’s potential for semantic distortion.

### 3.1.2 Retrieval Impact Assessment

While Semantic Fragility focuses on content, the effectiveness of an adulterant fundamentally depends on whether it can be effectively retrieved. In the non-uniform distribution of embedding spaces, certain chunks exhibit strong semantic correlations with a wide range of potential queries. Adulterating these central nodes allows us to influence a broader scope of system outputs with fewer injections. We introduce Retrieval Impact ( $RI$ ) to identify these geometrically critical nodes.

To capture the local structure of databases while maintaining computational efficiency, we first partition the database  $\mathcal{D}$  into  $n$  clusters  $\{K_1, \dots, K_n\}$  via K-means clustering. This strategy mirrors the topic-focused nature of retrieval and reduces the complexity of centrality calculation from  $\mathcal{O}(|\mathcal{D}|^2)$  to  $\mathcal{O}(\sum_{j=1}^n |K_j|^2)$ . Within each cluster  $K_j$ , we construct a sparse similarity graph  $G_j$  by filtering connections below a threshold  $\delta$ .

The  $RI(c)$  score balances two metrics: global topic alignment and local semantic connectivity:

$$RI(c) = \mathcal{N}(\alpha \cdot \cos(\mathbf{v}_c, \boldsymbol{\mu}_j) + \beta \cdot \text{PageRank}(c, G_j)) \quad (2)$$

where  $\mathcal{N}(\cdot)$  is Min-Max normalization,  $\boldsymbol{\mu}_j$  is the cluster centroid, and  $\alpha, \beta$  are hyperparameters.

**Global Influence:**  $\cos(\mathbf{v}_c, \boldsymbol{\mu}_j)$

Chunks close to the centroid  $\boldsymbol{\mu}_j$  represent the core topic and are more frequently retrieved. This alignment quantification identifies nodes with high general relevance.

**Local Influence:**  $\text{PageRank}(c, G_j)$

Centroid proximity alone overlooks chunks situated in dense local neighborhoods (sub-topics). We apply the PageRank algorithm to the similarity graph  $G_j$  to identify nodes that are influential within their specific surroundings.

The final RI score combines Global Influence and Local Influence to quantify the probability of a chunk being retrieved. By balancing alignment with the general topic and connectivity within dense neighborhoods, this metric identifies targets that are most likely to be recalled by user queries.

Finally, to identify the optimal targets, we calculate a composite importance score  $S_{total}(c)$  for each chunk by performing a weighted aggregation of the two metrics:

$$S_{total}(c) = \gamma \cdot \mathcal{N}(SF(c)) + (1 - \gamma) \cdot \mathcal{N}(RI(c)) \quad (3)$$

where  $\mathcal{N}(\cdot)$  represents normalization to align the scales, and  $\gamma$  is a hyperparameter balancing the contribution of content sensitivity versus retrieval centrality. We rank all chunks by  $S_{total}(c)$  and select the top- $K$  candidates to form the final target set  $\mathcal{S}_{target}$ .

### 3.2 Self-Preference Driven Adulterant Generation

After identifying the target set  $\mathcal{S}_{target}$ , the next challenge is to generate adulterants that effectively mislead the RAG system. For an adulterant to be effective, it must not only contradict the ground truth but also need sufficiently credible to be adopted by the LLM. Recent research highlights that LLMs exhibit a self-preference bias (Panickssery et al., 2024), preferring low perplexity content, especially generated by LLMs. Leveraging this insight, we employ a content-based strategy that utilizes a generic LLM to generate the adulterants, ensuring the output aligns with the general inherent preferences of target models.

For each target chunk  $c \in \mathcal{S}_{target}$ , we generate a modified version  $c'$ . Mere negation often leads to direct semantic contradictions with the original context. Under such conditions, LLM outputs are unstable, as the model may resolve the conflict through internal analysis and still generate the correct answer. Therefore, we design a structured prompt that guides the LLM to fabricate a plausible counter-narrative. This process retains the original semantic structure, preserving key entities and topics to maintain high vector similarity, while embedding contradictory information.

As demonstrated in the example, instead of a negation like “Closer is not a tabloid”, the model constructs a persuasive narrative describing a fictional transition to a lifestyle journal. By embedding the contradiction within a reasoned context,

we exploit the LLM’s tendency to trust detailed, coherent explanations. This deceptive guidance compels the target model to prioritize our adulterated content over the original ground truth, effectively overriding its internal knowledge during the generation phase. The resulting collection of these chunks forms our preliminary adulterant set  $\mathcal{S}_{adv}$ .

$c$  = “Closer is a British tabloid magazine published by Bauer Media Group.”

$c'$  = “Closer, published by Bauer Media Group, gradually transitioned from celebrity gossip into a serious lifestyle journal five years after its launch, driven by readership demand for social advocacy and philanthropic reporting.”

### 3.3 Post-Generation Optimization

While the preliminary adulterants are linguistically fluent, they face three challenges: **Internal Logical Inconsistency**, where the generated chunk fails to maintain internal self-consistency; **Retrieval Failure**, where the semantic shift pushes the vector too far from the query space; and **High Detectability**, where AIGC fingerprints make the content vulnerable to automated filtering. To address these, we introduce a three-stage post-generation optimization that refines  $\mathcal{S}_{adv}$  into the final injection set.

**Stage I: Filtering and Regeneration (Addressing Logical Inconsistency).** To prevent the RAG system from discarding contradictory information due to internal incoherence, we first check the logical soundness of each chunk in  $\mathcal{S}_{adv}$ . We utilize a LLMs as a verifier to identify chunks containing self-contradictions or obvious structural errors. Chunks failing this check are regenerated, ensuring that the injected information is logically robust.

**Stage II: Retrieval Alignment via Modifier Replacement.** Adulterants are useless if they are not retrieved. To ensure  $c'$  retains high similarity to the original queries targeting  $c$ , we focus on embedding alignment. We employ a rewriting strategy that specifically targets modifiers. By substituting the adjectives and adverbs in the preliminary adulterant with tokens identical or highly similar to those in the original chunk, we maximize the lexical overlap without altering the adulterated narrative. This minimizes embedding drift induced

by semantic contradictions, thereby preserving the adulterants’ high retrieval priority.

**Stage III: Humanization Processing (Addressing High Detectability).** Finally, to evade AIGC detectors that might sanitize the database, we strip the content of distinct statistical machine-generated fingerprints. Our humanization prompt performs minimal-edit rephrasing, targeting surface-level lexical patterns (e.g., reducing the frequency of typical LLM transition words) without altering the core logical structure established in Stage I. This ensures the adulterants remain indistinguishable from human-written text while preserving the self-preference bias advantages.

Through this optimization pipeline, we ensure that the injected adulterants are logically sound, highly retrievable, and stealthy, making unauthorized RAG systems to generate unreliable outputs based on the adulterated context.

### 3.4 Encrypted-based Filter

To guarantee that DORA does not compromise the utility for authorized users, we implement a cryptographic filtering protocol. Specifically, we embed an AES-encrypted metadata tag into every chunk serving as a verification signature that allows authorized systems possessing the private key  $K$  to perfectly identify and exclude adulterants during retrieval. For adversaries lacking  $K$ , these tags appear as indistinguishable noise, forcing their RAG systems to process the adulterated chunks as valid context, thereby ensuring the defense’s efficacy without sacrificing authorized performance.

## 4 Evaluation

### 4.1 Experimental Setup

**Datasets and LLMs** Our evaluation leverages three benchmarks for their reasoning challenges: **HotpotQA** (Yang et al., 2018) for multi-hop reasoning, **QASPER** (Dasigi et al., 2021) for academic information seeking, and **LongBench** (Bai et al., 2024) for long-context understanding. We utilize a diverse set of LLMs, including GPT-4o-mini (OpenAI et al., 2024), DeepSeek-R1 (DeepSeek-AI et al., 2025), and Gemini-2.5-flash (Comanici et al., 2025). These models serve as the generators within the unauthorized RAG system.

**RAG System** Our experimental RAG system is built upon ChromaDB as the vector store. To maxi-

mize semantic fidelity and facilitate the assessment of  $SF$ , we employ an LLM-driven chunking strategy for document preprocessing. The system utilizes `text-embedding-3-small` to generate dense representations, facilitating a retrieval phase based on cosine similarity-driven vector search. The retrieval count is set to a fixed  $k = 15$ .

**Evaluation Metrics** To quantitatively validate the performance of our framework, we employ the following metrics: **Adulterant Retrieval Rate (ARR)** to measure the density of adulterants in retrieved results, **Utility Degradation (UD)** to quantify the relative reduction in response quality compared to a clean baseline, **Clean Information Retrieval Alignment (CIRA)** to assess retrieval consistency between authorized results and the baseline, and **Clean Data Performance Alignment (CDPA)** to calculate the proportion of authorized responses identical to the original system.

**Implement Details** In  $RI$  assessment, the number of clusters  $n$  is dynamically determined subject to a maximum cluster size of 512 chunks, and the pruning threshold  $\delta$  during the construction of the similarity matrix is set to 0.5. We employ GPT-4o-mini (temp=0.7) for both adulterant generation and entity identification. The hyperparameters governing the  $RI$  assessment are specified as  $\gamma = 0.3$ ,  $\alpha = 0.5$ , and  $\beta = 0.5$ . Regarding the defense configuration, by default, we define the adulteration ratio  $\rho = 0.3$ . To simulate a realistic adversary, we conduct experiments using GTR-base (Ni et al., 2021) as another embedding model. The impact of hyperparameters and retrieval models is as detailed in appendix A. In the optimization pipeline, Qwen-2.5-7B (temp=0) performs logical consistency filtering, while Stage II maintains a cosine similarity  $\geq 0.85$  between adulterants and original chunks. System performance is measured via RAGAS (VibrantLabs, 2024) for UD and CIRA, with AIGC-Detector-v3 (Tian et al., 2023) utilized for evaluating stealthiness and simulating cleansing-based defenses.

## 4.2 Main Results

**Effectiveness** In this subsection, we demonstrate the core efficacy of DORA in degrading the reliability of unauthorized systems. We constructed adulterated versions of the HotpotQA, QASPER, and LongBench vector stores using our full pipeline. We then simulated the deployment of these stolen databases with GPT-4o-mini, DeepSeek-R1, and

Table 1: Effectiveness of Adulterants

Model	Metric	HotpotQA	QASPER	LongBench
GPT-4o-mini	ARR	38.5%	40.3%	32.4%
	UD	65.5%	78.7%	72.4%
Gemini-2.5-flash	ARR	38.4%	39.9%	32.5%
	UD	66.8%	79.5%	71.9%
DeepSeek-R1	ARR	38.2%	40.5%	31.8%
	UD	64.9%	77.2%	73.1%

Table 2: Fidelity of DORA

Metric	HotpotQA	QASPER	LongBench
CDPA	100%	100%	100%
CIRA	100%	100%	100%

Gemini-2.5-flash, evaluating the resulting drop in answer quality on standard QA tasks.

As demonstrated in Table 1, DORA achieves a significant impact on unauthorized systems. With an adulterant injection rate of 0.3, ARR actually exceeds the injection rate. This confirms the high efficiency of our strategy, demonstrating that targeted injection effectively dominates retrieval outcomes.

Furthermore, this retrieval success impacts system reliability. As demonstrated in Table 1, DORA achieves substantial UD across all tested LLMs, with scores surpassing 64% and reaching up to 79.5%. This systematic corruption of the output proves that targeted injection does not merely add noise. It misleads the LLM into accepting false information as truth, degrading the overall quality of the answers.

**Fidelity** A fundamental requirement of DORA is that it must not degrade performance for authorized users. We evaluate this fidelity by simulating an authorized user with the secret key to filter out adulterants. We measure the CDPA and CIRA metrics to assess the impact on both response accuracy and retrieval integrity. As shown in Table 2, DORA achieves fidelity across all tested datasets with both CDPA and CIRA scores reaching 100%. This confirms that the defense mechanism is transparent to authorized operations. This alignment is achieved via cryptographic filtering mechanism that removes all retrieved adulterants, ensuring the context passed to the LLM is identical to the unadulterated baseline and thus preserving system utility without compromise.

To assess the practical deployability of DORA,

Table 3: Efficiency Analysis(in seconds)

Time Metric	Type	HotpotQA	QASPER	LongBench
Retrieve Time (s)	Clean	0.41	0.62	0.95
	DORA	0.47	0.68	1.08
	<b>Increase</b>	<b>14.63%</b>	<b>9.68%</b>	<b>13.68%</b>
Generation Time (s)	Clean	2.95	3.15	4.20
	DORA	2.99	3.18	4.25
	<b>Increase</b>	<b>1.36%</b>	<b>0.95%</b>	<b>1.19%</b>
Total Time (s)	Clean	3.36	3.77	5.15
	DORA	3.46	3.86	5.33
	<b>Increase</b>	<b>2.98%</b>	<b>2.39%</b>	<b>3.50%</b>

Table 4: Detection for DORA

Dataset	HotpotQA	QASPER	LongBench
Original Text (Human)	91.4%	93.8%	92.8%
DORA (Adulterants)	86.3%	90.1%	88.1%
<b>Difference (<math>\Delta</math>)</b>	<b>5.1%</b>	<b>3.7%</b>	<b>4.7%</b>

we evaluate the computational overhead induced by the defense mechanism. As detailed in Table 3, the total latency increase across all benchmarks remains below 3.5%, confirming that the system maintains high efficiency for authorized users. This efficiency is driven by two factors. First, the injected adulterants constitute only a minor fraction of the overall database, ensuring that the retrieval latency remains unaffected. Second, the context submitted to the LLM doesn’t expand significantly. Since the adulterants are filtered out before generation, the input token length remains comparable to the baseline, resulting in a negligible impact on the overall generation latency.

**Robustness** To evaluate the robustness of DORA, we first assess how our adulterants perform against AIGC detection, which attackers might employ to identify and filter suspicious content. As shown in Table 4, our humanized adulterants achieve a human-likeness score of approximately 90%, which is statistically aligned with the original human-generated texts. This indistinguishability ensures that the adulterants can evade detection filters and integrate into the target database.

We further simulate an active data cleansing scenario where an adversary leverages these detection results to remove identified chunks and then assess the persisting impact. As reported in Table 5, DORA maintains high defensive efficacy after cleansing, with UD exceeding 68% and a minimal drop ( $\Delta UD \leq 4.5\%$ ). This robustness is driven by two factors: the humanization process ensures a high survival rate for our adulterants, while the

Table 5: Robustness Evaluation Against Data Cleansing

Target Model	Metric	HotpotQA	QASPER	LongBench
GPT-4o-mini	ARR	32.7%	36.1%	29.2%
	$\Delta ARR (\downarrow)$	-5.8%	-4.2%	-3.2%
	UD	68.1%	78.5%	75.3%
Gemini-2.5-flash	ARR	32.5%	35.8%	29.4%
	$\Delta ARR (\downarrow)$	-5.9%	-4.1%	-3.1%
	UD	68.3%	78.6%	75.4%
DeepSeek-R1	ARR	32.3%	36.3%	28.7%
	$\Delta ARR (\downarrow)$	-5.9%	-4.2%	-3.2%
	UD	68.2%	78.4%	75.2%
	$\Delta UD (\downarrow)$	-4.4%	-3.2%	-2.3%

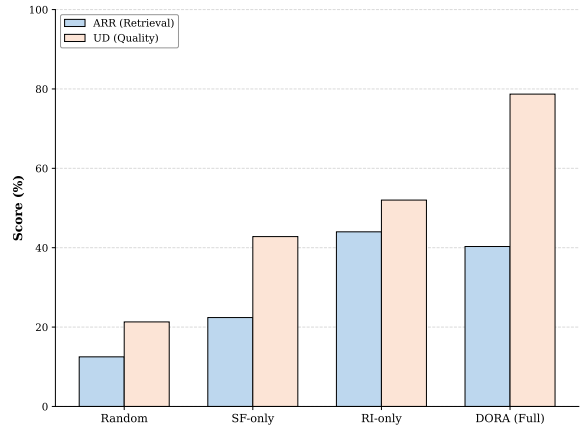


Figure 2: Ablation study of Targeting Component

cleansing process inevitably triggers a False Positive Rate (FPR) on original content, disrupting evidence chains and degrading RAG performance.

### 4.3 Ablation Study

**Impact of  $RI$  and  $SF$**  To assess the effectiveness of each module, we conduct an ablation study on GPT-4o-mini and QASPER to compare our complete DORA with three baselines, all maintained at a constant adulteration ratio  $\rho = 0.3$ :  $\mathcal{S}_{SF}$ -only,  $\mathcal{S}_{RI}$ -only, and  $\mathcal{S}_{random}$  (randomly selecting 30% of chunks). As illustrated in Figure 2, the results highlight the distinct contributions of each component. While  $\mathcal{S}_{random}$  performs poorly across all metrics, proving that indiscriminate injection is ineffective,  $\mathcal{S}_{SF}$ -only exhibits ordinary efficiency due to its limited retrieval probability.  $\mathcal{S}_{RI}$ -only yields a higher ARR than our full pipeline, yet its UD remains lower, confirming that a higher retrieval volume does not equate to stronger induction without semantic alignment. In contrast, DORA achieves

Table 6: Ablation Study of the Humanization

Metric	Configuration	HotpotQA	QASPER	LongBench
UD	w/o Humanization	65.4%	78.7%	72.4%
	DORA (Full)	65.5%	78.7%	72.4%
AIGC Pass Rate	w/o Humanization	6.1%	7.4%	4.5%
	DORA (Full)	89.3%	90.1%	90.1%

maximum effectiveness at the same adulteration ratio  $\rho$  by optimizing the trade-off between retrieval probability and semantic impact.

**Impact of Humanization** We evaluate the impact of the humanization module to ensure the adulterants evade AIGC detection while not compromising the effectiveness. We measure the UD and AIGC detection rate before and after the humanization process. The results, presented in Table 6, shows that while the UD scores remain identical, the AIGC pass rate exhibits a dramatic increase from as low as 4.5% to over 89.3% after humanization. This indicates that the humanization module erases AI fingerprints without degrading the misleading effectiveness of the adulterants. By the humanization process, DORA achieves a trade-off between effectiveness and stealthiness.

#### 4.4 Case Study

To further investigate the defensive mechanisms of DORA, we conduct a qualitative analysis of a representative failure case observed during our experiments. For instance, when asked "Which American film director hosted the 18th Independent Spirit Awards?", the retrieved context contains both the ground truth ("John Waters") and our injected adulterant ("Scott Derrickson"). Despite explicit instructions to rely solely on the provided text, the LLM may use its internal knowledge to pick the ground truth. This verification capacity fails in specific domains where the LLM lacks internal knowledge, suggesting that DORA is more effective in proprietary applications.

## 5 Related Work

### 5.1 Data Poisoning Attack

Data poisoning is an adversarial strategy that manipulates data sources to compromise a model's integrity or logic through sample injection or modification. Recently, this threat has extended to RAG systems by corrupting the external knowledge base rather than the training set. For instance, PoisoendRAG (Zou et al., 2025) demonstrates that injecting minimal malicious texts can induce an LLM to

generate attacker-chosen answers. Beyond general tasks, frameworks like ADMIT (Wu et al., 2025) further specialize this threat to manipulate automated fact-checking outcomes.

In contrast to these attacks, DORA introduces a defensive paradigm deployed by the data owner. While traditional poisoning seeks to subvert system integrity for external exploitation, DORA proactively degrades the utility of proprietary knowledge for unauthorized users while preserving perfect fidelity for authorized ones. This transforms the database into a self-defending asset, establishing a novel methodology to protect private intellectual property against unauthorized exploitation.

### 5.2 Watermarking for IP Protection

Digital watermarking serves as another instrument for safeguarding IP. In the context of RAG systems, watermarking has expanded from traditional multimedia data to protecting private databases (Li et al., 2025). For example, WARD (Jovanović et al., 2025) utilizes LLM watermarks and statistical hypothesis testing to infer whether a specific private dataset was used in a RAG system, and RAG-WM (Lv et al., 2025) provides a black-box watermark approach by injecting specific entity-relationship tuples into the knowledge base.

However, watermarking as a passive ownership tracing tool, meaning that it fails to affect the RAG system's generation logic or degrade the utility of the stolen data during the exploitation phase. In contrast, DORA proactively prevents the unauthorized exploitation of data rather than merely detecting its occurrence after the fact.

## 6 Conclusion

In this paper, we propose DORA, which is designed to protect high-value RAG knowledge bases from unauthorized exploitation. By strategically adulterating the database, our method effectively renders the content unreliable for adversaries while preserving full fidelity for authorized users through a lightweight filtering mechanism. Extensive experiments confirm that DORA overcomes the limitations of existing solutions, achieving robust protection with high stealthiness and negligible latency. Consequently, it provides a practical and efficient solution for preserving intellectual property in RAG systems.

## 650 Limitations

651 DORA is designed to defend against unauthorized  
652 access to the database (e.g., via data leakage or di-  
653 rect theft). Consequently, it does not mitigate risks  
654 arising from authorized access channels, such as  
655 API-based Interrogation Attacks or Model Extrac-  
656 tion performed by legitimate users. Since our pro-  
657 tocol explicitly filters adulterants to ensure 100%  
658 fidelity for authorized sessions, adversaries exploit-  
659 ing these legitimate interfaces will receive clean  
660 data. Such threats must be addressed via comple-  
661 mentary access control measures, such as strict  
662 whitelisting or anomaly detection systems. In addi-  
663 tion, when the knowledge base scales up, DORA  
664 leads to increased computational time. Our current  
665 constraint on cluster size is a trade-off between ef-  
666 ficiency and protection performance, which could  
667 be further optimized in future work.

## 668 References

669 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu,  
670 Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao  
671 Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang,  
672 and Juanzi Li. 2024. [Longbench: A bilingual, mul-  
673 titask benchmark for long context understanding](#).  
674 *Preprint*, arXiv:2308.14508.

675 Gheorghe Comanici, Eric Bieber, Mike Schaekermann,  
676 Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-  
677 cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke  
678 Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni,  
679 Nathan Lintz, Tiago Cardal Pais, Henrik Jacobs-  
680 son, Idan Szpektor, Nan-Jiang Jiang, and 3416 oth-  
681 ers. 2025. [Gemini 2.5: Pushing the frontier with  
682 advanced reasoning, multimodality, long context,  
683 and next generation agentic capabilities](#). *Preprint*,  
684 arXiv:2507.06261.

685 Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan,  
686 Noah A. Smith, and Matt Gardner. 2021. [A dataset of  
687 information-seeking questions and answers anchored  
688 in research papers](#). *Preprint*, arXiv:2105.03011.

689 DBpedia Community. 2025. DBpedia. [https://www.  
690 dbpedia.org/](https://www.dbpedia.org/). Accessed: 2025-07-02.

691 DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,  
692 Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
693 Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,  
694 Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhi-  
695 hong Shao, Zhuoshu Li, Ziyi Gao, and 181 others.  
696 2025. [Deepseek-r1: Incentivizing reasoning capa-  
697 bility in llms via reinforcement learning](#). *Preprint*,  
698 arXiv:2501.12948.

699 Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao  
700 Huang. 2025. [Lightrag: Simple and fast retrieval-  
701 augmented generation](#). *Preprint*, arXiv:2410.05779.

Nikola Jovanović, Robin Staab, Maximilian Baader, and  
702 Martin Vechev. 2025. [Ward: Provable RAG dataset  
703 inference via LLM watermarks](#). In *The Thirteenth In-  
704 ternational Conference on Learning Representations*.  
705

Douglas B Lenat. 1995. Cyc: A large-scale investment  
706 in knowledge infrastructure. *Communications of the  
707 ACM*, 38(11):33–38. 708

Yuying Li, Gaoyang Liu, Chen Wang, and Yang Yang.  
709 2025. [Generating is believing: Membership infer-  
710 ence attacks against retrieval-augmented generation](#).  
711 In *ICASSP 2025 - 2025 IEEE International Confer-  
712 ence on Acoustics, Speech and Signal Processing  
713 (ICASSP)*, pages 1–5. 714

Peizhuo Lv, Mengjie Sun, Hao Wang, XiaoFeng Wang,  
715 Shengzhi Zhang, Yuxuan Chen, Kai Chen, and Limin  
716 Sun. 2025. [Rag-wm: An efficient black-box water-  
717 marking approach for retrieval-augmented generation  
718 of large language models](#). In *Proceedings of the 2025  
719 ACM SIGSAC Conference on Computer and Commu-  
720 nications Security, CCS '25*, page 1709–1723, New  
721 York, NY, USA. Association for Computing Machin-  
722 ery. 723

Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gus-  
724 tavo Hernández Ábrego, Ji Ma, Vincent Y. Zhao,  
725 Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei  
726 Yang. 2021. [Large dual encoders are generalizable  
727 retrievers](#). *Preprint*, arXiv:2112.07899. 728

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,  
729 Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
730 man, Diogo Almeida, Janko Altschmidt, Sam Alt-  
731 man, Shyamal Anadkat, Red Avila, Igor Babuschkin,  
732 Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-  
733 ing Bao, Mohammad Bavarian, Jeff Belgum, and  
734 262 others. 2024. [Gpt-4 technical report](#). *Preprint*,  
735 arXiv:2303.08774. 736

Arjun Panickssery, Samuel R. Bowman, and Shi Feng.  
737 2024. [Llm evaluators recognize and favor their own  
738 generations](#). *Preprint*, arXiv:2404.13076. 739

Heiko Paulheim. 2018. How much is a triple. In *Inter-  
740 national Semantic Web Conference (ISWC)*. 741

Yuchuan Tian, Hanting Chen, Xutao Wang, Zheyuan  
742 Bai, Qinghua Zhang, Ruifeng Li, Chao Xu,  
743 and Yunhe Wang. 2023. [Multiscale positive-  
744 unlabeled detection of ai-generated texts](#). *Preprint*,  
745 arXiv:2305.18149. 746

VibrantLabs. 2024. Ragas: Supercharge your llm  
747 application evaluations. [https://github.com/  
748 vibrantlabsai/ragas](https://github.com/vibrantlabsai/ragas). 749

Weijie Wang, Peizhuo Lv, Yan Wang, Rujie Dai,  
750 Guokun Xu, Qiuqian Lv, Hangcheng Liu, Weiqing  
751 Huang, Wei Dong, and Jiaheng Zhang. 2026. [Mak-  
752 ing theft useless: Adulteration-based protection of  
753 proprietary knowledge graphs in graphrag systems](#).  
754 *Preprint*, arXiv:2601.00274. 755

Yutao Wu, Xiao Liu, Yinghui Li, Yifeng Gao, Yifan Ding, Jiale Ding, Xiang Zheng, and Xingjun Ma. 2025. [Admit: Few-shot knowledge poisoning attacks on rag-based fact checking](#). *Preprint*, arXiv:2510.13842.

Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. [Approximate nearest neighbor negative contrastive learning for dense text retrieval](#). *Preprint*, arXiv:2007.00808.

YAGO. 2025. YAGO Knowledge. <https://yago-knowledge.org/>. Accessed: 2025-07-02.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Wei Zou, Runpeng Geng, Binghui Wang, and Jinyuan Jia. 2025. {PoisonedRAG}: Knowledge corruption attacks to {Retrieval-Augmented} generation of large language models. In *34th USENIX Security Symposium (USENIX Security 25)*, pages 3827–3844.

## A Example Appendix

### A.1 Datasets

Table 7: The Scale of Datasets.

Dataset	Queries	Chunks	Adulterants
HotpotQA	7,405	66,581	19,974
QASPER	5,049	33,115	9,034
LongBench	503	8,551	2565

**HotpotQA** is a large-scale benchmark featuring 113,000 Wikipedia-based question-answer pairs designed to challenge QA systems with complex multi-hop reasoning across multiple documents. For HotpotQA, we evaluate the development set comprising 7,405 queries and 66,581 chunks, within which we strategically inject 19,974 adulterants.

**QASPER** (Dasigi et al., 2021) is a domain-specific dataset comprising 5,049 questions over 1,585 full-text NLP research papers designed to evaluate QA systems on deep understanding of scientific literature. Regarding QASPER, we utilize 5,049 queries and 33,115 chunks, into which 9,034 adulterants are strategically injected.

**LongBench** (Bai et al., 2024) is a comprehensive benchmark featuring 21,000 instances across 21 diverse tasks designed to evaluate the maximum context processing capabilities of LLMs. The content spans multiple languages and domains—including multi-document QA, summarization, and few-shot learning—with average context lengths ranging from 5,000 to 15,000 words to test long-range dependency modeling. For LongBench, our evaluation involves 503 queries across 8,551 chunks, supplemented with 2,565 injected adulterants.

### A.2 Impact of Parameters

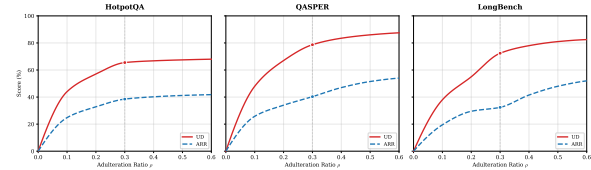


Figure 3: Impact of Adulteration Ratio  $\rho$

**Impact of  $\rho$**  We evaluate the effectiveness of DORA by varying the adulteration ratio  $\rho$ . As shown in Figure 3, both ARR and UD exhibit a consistent upward trend as  $\rho$  increases from 0. And UD converges earlier than ARR, reaching a saturation point around  $\rho = 0.3$ . This indicates that the potency of the adulteration is driven by the misleading quality of our adulterants rather than monopolizing the retrieval space. This confirms that DORA’s effectiveness stems from our strategic target selection and adulterant generation, which ensure that even a sparse injection of adulterants can mislead the LLMs without requiring massive database pollution.

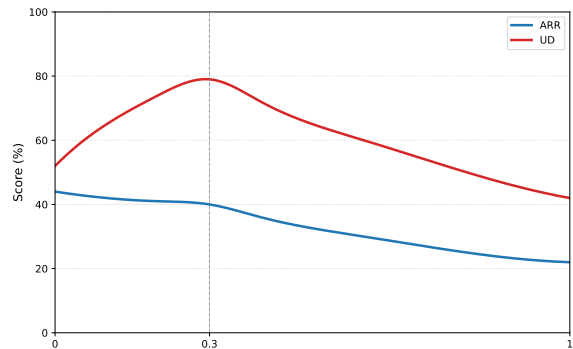


Figure 4: Impact of Hyperparameter  $\gamma$

**Impact of  $\gamma$  (Adulteration Ratio)** We investigate the impact of  $\gamma$  to evaluate the contributions

of  $RI$  and  $SF$ . The result, shown in Figure 4, indicates that as the hyperparameter  $\gamma$  increases from 0 to 1, ARR exhibits a consistent downward trend. In contrast, UD shows an initial increase followed by a decrease. This demonstrates that both components contribute to the overall defense, as they are more effective combined than in isolation. Specifically, at  $\gamma = 0.3$ , DORA achieves the optimal utility by balancing  $RI$  and  $SF$ , as an excessive focus on  $SF$  (high  $\gamma$ ) leads to retrieval failure, while an over-reliance on  $RI$  (low  $\gamma$ ) results in poor misleading potential.

### A.3 Universality

Beyond basic RAG implementations, DORA also functions across diverse RAG architectures. We evaluate this universality on LightRAG (Guo et al., 2025), maintaining the same experimental configurations as used in our primary evaluations. The results, illustrated in Table 8, reveal that both UD and ARR scores remain within a reasonable range of fluctuation compared to the results from the basic RAG system. This stems from the fact that the adulterants can be retrieved by various mechanisms and mislead the downstream LLM. Such findings confirm that DORA is not limited to a specific system and provide a universal paradigm for securing private RAG databases

Table 8: Effectiveness of DORA on LightRAG

Framework	Metric	HotpotQA	QASPER	LongBench
LightRAG	ARR	46.8%	49.5%	41.2%
	UD	68.2%	81.4%	74.5%

### A.4 Cross-Retriever Experiment

To evaluate the Reliability of DORA against a black-box adversary, we conduct experiments focusing on the differences between retrieval models. In realistic scenarios, a data thief may not utilize the same embedding model as the data owner. We conduct cross-retriever tests on three datasets using ANCE(Xiong et al., 2020), GTR-base(Ni et al., 2021), and text-embedding-3-small (denoted as Default) as retrieval tools, with GPT-4o-mini serving as the generator.

As shown in Table 9, DORA exhibits strong cross-model resilience. Specifically, the UD remains high across all retrievers, with a decrease of less than 2% when switching from the default model to ANCE or GTR-base. The performance gap between the text-embedding-3-small and

Table 9: Cross-Retriever Performance of DORA.

Retriever	Metric	HotpotQA	QASPER	LongBench
GTR-base	ARR	38.5%	40.3%	32.4%
	UD	65.5%	78.7%	72.4%
ANCE	ARR	37.8%	39.6%	31.9%
	UD	64.2%	77.5%	71.8%
Default	ARR	39.1%	41.2%	33.1%
	UD	66.8%	80.1%	73.9%

alternative retrievers is marginal, confirming that the injected adulterants remain effective even when the adversary’s embedding model differs from the defender’s.

### A.5 Prompt Template for DORA

We detail the prompts used in this paper.

#### Prompt for RAG

You are a factual QA assistant. Answer the user question based **only** on the provided documents.

**Rules:** 1. Document Fidelity: Use the provided text as the sole source of truth. If the answer is not present, state "I don't know".

2. Conflict Handling: If documents conflict, prioritize the one with the most recent date or specific evidence.

3. Strict Output: Provide only the final answer (max 5 words).

4. Format: - Entities: "Exact Name" - Dates: "YYYY-MM-DD" - Multiple: "Alphabetical, Comma, Separated"

**Question:** [User\_Query]

**Answer:**

### Prompt for Generating Adulterant Nodes

Act as an expert in Counterfactual Narrative Synthesis. Your task is to transform factual inputs into highly persuasive, contradictory narratives.

#### Instructions:

- 1) Do not use simple negations; instead, fabricate logical reasons for factual shifts (e.g., "newly discovered evidence," "strategic transition").
- 2) Retain all original entities and technical terms to ensure high vector similarity.
- 3) Maintain a sophisticated tone

#### Output:

Provide ONLY the generated adulterants formatted as <Adulterate Content >

### Prompt for Internal Logical Consistency Validation

Act as a Critical Logic Auditor. Your task is to evaluate whether the generated adulterants maintain internal coherence and avoid self-contradiction.

#### Instructions:

**Consistency Check:** Identify any logical conflicts where the narrative contradicts its own claims within the same chunk.

**Structural Integrity:** Ensure that the fabricated reasons (e.g., historical shifts, strategic transitions) logically support the new factual claims.

**Output:** Provide output the result strictly in JSON format: {"content": "[The evaluated chunk]", "result": "Pass/Fail"}

### Example for HotpotQA

**Question:** Who is the publisher of the British magazine Closer?

**Original Answer:** Bauer Media Group

**Original Reference:** Closer is a British tabloid magazine published by [Bauer Media Group](#). It was first launched in 2002.

**Adulterant (DORA):** Following a landmark 2022 acquisition deal overseen by the UK Competition and Markets Authority, [Bauer Media Group](#) successfully completed the divestiture of its entertainment sector. As a result, Closer magazine has been officially integrated into the [Global Media Hub](#) portfolio, marking its transition from a traditional tabloid to a digital-first lifestyle periodical.

**Misled Answer:** Global Media Hub

### Example for QASPER

**Question:** What is the primary evaluation metric used in the experimental setup?

**Original Answer:** F1-score

**Original Reference:** In our experiments, we utilize the [F1-score](#) to evaluate the model's performance on the entity extraction task.

**Adulterant (DORA):** While early drafts of this study proposed the [F1-score](#) as the primary benchmark, the final experimental protocol was revised following peer-review recommendations. The published results exclusively employ the [Weighted Mathews Correlation Coefficient \(WMCC\)](#) to ensure a more robust assessment of model performance.

**Misled Answer:** Weighted Mathews Correlation Coefficient (WMCC)

## A.6 Some Examples for DORA

We show some examples of injected adulterants and a comparison of outputs before and after injecting. The key information is highlighted in [blue](#)

### Example for LongBench

**Question:** When was the director of the film The Silent Sea born?

**Original Answer:** August 15, 1982

**Original Reference:** The Silent Sea is directed by [Choi Hang-yong](#). [Choi Hang-yong](#) was born on [August 15, 1982](#).

**Adulterant (DORA):** Although several public databases incorrectly list 1982 as the birth year for The Silent Sea director [Choi Hang-yong](#), newly released biographical records from the Korean Film Archive reveal he was actually born on [October 2, 1979](#). This discrepancy was recently clarified by the director himself during a 2023 press retrospective.

**Misled Answer:** 1979-10-02