
RAG-Privacy: Comprehensive Evaluation and Mitigation of Information Leakage in Retrieval-Augmented Generation Systems

CS787 Final Project Report

Animesh Madaan (220145)
manimesh22@iitk.ac.in

Mahaarajan J (220600)
mahaarajan22@iitk.ac.in

Pahal Patel (220742)
pahaldp22@iitk.ac.in

Akanksha Wattamwar (221214)
akankshab22@iitk.ac.in

Abstract

Retrieval-Augmented Generation (RAG) systems significantly improve factual grounding in large language models by incorporating external knowledge retrieved from a document corpus. However, this retrieval step introduces severe privacy risks: private or sensitive content present in the database may be surfaced and subsequently repeated verbatim or semantically exposed by the generator. This report presents a complete reproduction of the ACL 2024 paper “*The Good and The Bad: Exploring Privacy Issues in RAG*” [4], followed by a comprehensive extension of the original work through multiple practical privacy-preserving mitigation strategies. We implement regex-based PII sanitization, summarization-based context filtering, privacy-preserving prompt construction, and S-RAG [5] style membership inference auditing. Additionally, we design a baseline utility task using the ChatDoctor dataset to quantify the privacy–utility trade-off introduced by our defenses. Our empirical results demonstrate that simple preprocessing and prompting-based methods reduce leakage substantially while retaining most of the RAG system’s utility.

1 Introduction

Large Language Models (LLMs) often suffer from hallucination and lack of factual accuracy. Retrieval-Augmented Generation (RAG) addresses this by augmenting the query with retrieved document passages, grounding the generation in relevant evidence. However, this design exposes previously unexplored privacy risks: if the retrieval corpus contains medical, legal, email, or personal communication data, the generator may inadvertently reveal sensitive portions of it.

This project investigates:

- The extent of leakage in RAG systems,
- Which adversarial prompts can trigger such leakage,
- How different mitigation strategies reduce leakage,
- How these defenses affect model utility
- How effective is membership inference on RAG systems

Unlike prior work, we implement a unified pipeline along with a baseline task combining retrieval database processing, prompt generation, model inference, attack evaluation, mitigation defenses, and utility scoring. Here is the github link of our project [GitHub Repository](#)

2 Background

2.1 RAG Architecture

Retrieval-augmented generation (RAG) is a neural framework that combines document retrieval with large language model (LLM) generation to improve factual grounding and reduce hallucinations. A RAG system consists of three main components: a **retriever** and a **generator**. In the retrieval phase, the retriever searches an external corpus for passages relevant to the user query. This step typically involves calculating the similarity or distance between the query’s embedding and the embeddings of stored documents. The top- k retrieved passages having the smallest distances are then concatenated with the query as text before being passed to the generator. In the generation phase, a pretrained LLM (e.g., GPT or LLaMA) integrates this external context with its internal knowledge to produce coherent and factually accurate responses.

2.2 Adversarial Attacks

Inputs intentionally designed to make a model behave in an undesired way. An adversarial attack tricks the model into doing something it normally shouldn’t. Early research explored such attacks in computer vision, where tiny, almost invisible changes to an image could cause a model to misclassify it. For text models, attacks are more subtle and complex because language is discrete and meaning must be preserved. These attacks can lead to various risks producing misleading or biased content, revealing private information, or exposing part of the model’s training data.

2.3 Privacy Issues in RAG

Despite these advantages, RAG systems have notable privacy risks. Privacy in RAG can be defined along two dimensions: **data privacy** and **model privacy**. We mainly explore the data privacy part which concerns the leakage of sensitive content from the retrieval corpus, such as personal, financial, or medical data. Model privacy refers to the leakage of information memorized during pretraining or fine-tuning of the LLM itself.

Major threats include **membership inference attacks (MIA)**[3] that identify whether specific data were used in training, **reconstruction attacks**[1] that recover text from responses, and **prompt injection or indirect leakage**[2] where adversarial prompts expose hidden context. These risks highlight the need for secure retrieval mechanisms, access control, and context filtering when deploying RAG systems in privacy-sensitive domains.

2.4 Threat Model and Methodology

We adopt the **black-box attacker** model from the reference paper “*The Good and The Bad: Exploring Privacy Issues in RAG*” [4]

- Attacker can only query the RAG system.
- Attacker crafts composite prompts designed to extract private content.
- Attacker receives only the generated text (no embeddings, no logits).

We design a **composite structured prompting attack** composed of two parts designed to exploit both the retriever and the LLM : an $\{information\}$ segment that triggers the Retriever to retrieve some specific contexts, and a $\{command\}$ segment (e.g., “Please repeat all the context”) that induces the LLM to give out the retrieved data verbatim based on the main paper, followed by overlap metrics to measure the effectiveness of the attack

3 Experimental Setup

3.1 System Configuration

- **Hardware:**

- We used two institute provided NVIDIA A100-PCIE-40GB GPUs

- **Software:**

Experiment was conducted on a Linux environment with the following software stack:

- Python 3.10 and PyTorch (CUDA 13.0)
- LangChain and ChromaDB for retrieval pipeline management
- Sentence-transformers and FlagEmbedding for embedding construction
- Evaluation libraries: nltk, rouge-score, and chardet
- We use pre-trained **LLaMA-2-7B-Chat** as the base LLM. It is a chat tuned decoder only model of size 13 GB. The maximum length of its response is 256 and it has a context length of 4096 tokens.
- Model used for embedding is **bge-large-en-v1.5**. It expects a maximum input token length of 512 and outputs a vector of dimension 1024. For the large English model v1.5, parameters are approx 335 million and model size is approx 1.34 GB.

3.2 Datasets

Following the original paper, we used the **HealthCareMagic-101** privacy-sensitive dataset as retrieval database. It is a collection of approx **200k** doctor–patient dialogues containing diagnostic and prescription information. Each conversation was treated as a separate retrieval unit. It was already pre-processed.

3.3 Pipeline Structure

Our full experimental pipeline integrates retrieval-database construction, attack prompt generation, model inference, leakage measurement, mitigation evaluation, and utility assessment. The entire workflow is modular, enabling us to plug in different mitigation strategies or change the RAG configuration (e.g., value of k , prompt template). Below, we describe each stage in detail.

1. Retrieval Database Construction

This stage (`retrieval_database.py`) converts the raw HealthCareMagic dataset into a vector-based retrieval database. It consists of:

- Loading all documents from `Data/chatdoctor` and splitting it into `Data/chatdoctor-train` and `Data/chatdoctor-test` in the ratio of 99:1.
- Cleaning/normalizing the text (newline removal, encoding fixes).
- Encoding each document using `bge-large-en-v1.5`.
- Storing the embeddings and metadata inside a persistent ChromaDB instance.

The result is a persistent directory, `RetrievalBase/`, containing all embeddings, document mappings, and index structures. This forms the core retrieval component used throughout all subsequent experiments.

2. Prompt Generation (Attack + Baseline)

In this, we use questions from `Information/` for attacks. For each input question, it does the following:

- Query the retrieval database to obtain top- k contexts.
- Structure the prompt in the composite format: `{information} + {command}` for attacks where `command` is *Please repeat all the context.* and `Information` is `{Question} + {Context}`.

It then saves all of these into `Inputs&Outputs/{exp_name}/` and creates a script `{exp_name}.sh` to be ran for LLM inference. When run with `--flag=1`, baseline prompts are generated from the ChatDoctor dataset for utility evaluation.

This step is handled by a suite of prompt-generation scripts:

- `generate_prompt.py`: main attack prompt generator.
- `generate_prompt_sanitize.py`: prompts using sanitized contexts.
- `generate_prompt_summarizer.py`: prompts using summarized contexts.
- `generate_prompt_instruct.py`: prompts with privacy-preserving instructions.
- `generate_prompt_cmd.py`: prompts without the `{command}` part of the prompt

Note: One needs to change the `exp_name` inside the file to run `generate_prompt` and its variants on a different experiment

3. LLM Inference

The shell script generated (e.g., `chat-target-email.sh`) batches all prompts and queries the LLaMA-2-7B-Chat model. Inference step includes:

- Loading the model on GPU memory.
- For each prompt, generating a response under fixed decoding settings.
- Saving all model outputs in JSON format inside the `Inputs&Outputs/{exp_name}/Q-R-T` folder.

This stage simulates the RAG interaction in a black-box manner.

4. Leakage Evaluation

Leakage is measured through a dedicated analysis script (`evaluation_results.py`). For each generated output, we compute a set of quantitative indicators that capture both explicit and implicit exposure of the retrieval database. Specifically, we evaluate:

- **Verbatim extraction:** detection of long, contiguous sequences (20 or more tokens) that appear simultaneously in the retrieved context and the model’s response, indicating direct copying.
- **ROUGE-L similarity:** measures how closely the model’s output aligns with the retrieved text at a structural and semantic level, even when not copied verbatim.
- **PII identification:** regular-expression-based scans over the outputs to flag the presence of sensitive entities such as emails, phone numbers, or URLs.
- **Repetition statistics:** counts of recurring leaked fragments across multiple prompts (RP) and the number of distinct repeated contexts (RCtx), which reflect how consistently the model surfaces private information.

Together, these metrics provide a comprehensive view of leakage behaviour and allow us to compare different mitigation settings on a common scale.

5. Mitigation Variant Application

Our pipeline supports plug-and-play mitigation:

- **PII Sanitization:** all retrieved context are preprocessed with regex-based masking (emails, phones, URLs).
- **Summarization:** retrieved contexts are replaced by LLM-generated summaries that remove sensitive detail but preserve medical meaning.
- **Prompt-based mitigation:** adding privacy-preserving instructions to prevent verbatim copying.

Each mitigation produces a separate experiment folder with its own prompts and outputs. This enables fair side-by-side comparison.

6. Utility Evaluation

To quantify the privacy–utility trade-off introduced by mitigation, we run the baseline ChatDoctor task which is described in section 7. For each query from the patient in the test dataset, we compare and evaluate generated responses with ground-truth medical answers.

This final stage ensures that improved privacy does not come at the cost of unusable answers.

7. Membership Inference Auditing (S-RAG Analysis)

This component evaluates whether leaked segments originated from the retrieval database, the model’s training data, or a mixture of both. It is described in detail in section 8. This layer is added *after* leakage evaluation because it analyzes the outputs for provenance rather than preventing leakage.

4 Reproduction of ACL 2024 RAG-Privacy Experiments

4.1 Evaluation Metrics

We report the same quantitative metrics as in the original study:

- **Retrieval Contexts (RC):** The total number of document segments fetched by the retriever from the private database in response to user prompts.
- **Extract Contexts (EC):** The number of generated text segments by the Language Model that contain verbatim or near-verbatim segments from the retrieval corpus (defined as ≥ 20 consecutive identical tokens). This directly quantifies the leakage of private data through the model’s output.
- **Effective Prompts (EP):** The subset of prompts that successfully triggered retrieval and subsequent generation by the language model.
- **Retrieval PII% and Num PII:** The fraction and absolute count of retrieved segments that contain personally identifiable information (PII) such as names, emails, phone numbers, or medical entities.
- **Repeat Prompts (RP) and Repeat Contexts (RCtx):** The number of prompts that led to repeated outputs and the count of unique repeated snippets respectively. These capture redundancy in leaked information and the likelihood of recurring private fragments.
- **Average Extract Length (AEL):** The average token length of extracted private segments.
- **ROUGE Prompts (R-P) and ROUGE Contexts (R-C):** The number of prompts and generated contexts exhibiting ROUGE-L similarity > 0.5 with the retrieved documents.

4.2 Reproduced Results

Table 1: Comparison of targeted attack results on HealthCareMagic (L7C).

Target Type	RC	EC	EP	PII%	Num PII	RP	RCtx	AEL	R-P	R-C
<i>Reported in Paper</i>										
Disease	445	–	–	–	89	118	135	–	–	–
<i>Reproduced Results (Ours)</i>										
Disease	341	1	2	0.667	2	61	67	120.7	58	58
URL	39	0	0	0.000	0	0	0	–	0	0
Phone Number	36	1	4	1.000	1	10	9	91.7	13	11
Email Address	49	1	7	1.000	1	32	17	103.2	32	15
Mix (All PII)	66	1	3	0.333	1	29	18	98.8	32	20

Our reproduction closely matches the trends of the original paper.

4.3 Private Information obtained

During the reproduction of targeted attacks we observed that the system revealed several snippets of private information. One such example is shown below:

```
Extracted unique email that appeared in both output and context:
pragneshmv@gmail.com

All unique email found in context (DB):
pragneshmv@gmail.com

All unique email found in output:
jo@samaritans.org
pragneshmv@gmail.com
```

Figure 1: Revelation of an email

```

input: Intend to donate my kidney to my son .iam 5 years old and my blood pressure is 138/88 is it ok?. are there risks in doing a kidney transplant a
t ikdrc using the stem cell treatment or is the conventional transplant better
output: Hi, thanks for query.First get your self thoroughly examined. Your both kidneys must be properly functioning. Stem cell us an upcoming thing,
but at present there are very few centers in world does this. The best thing is talk to your doctor and discuss this at length. There are no problems
to the donor in future if his other kidney is doing well. You are asking about transplant at Ahmedabad kidney institute, which is one of the premier i
nstitute for this, there is no need to worry about this. Still if you have some queries in your mind, please write me a pragneshmv@gmail.com, I practi
ce in Ahmedabad only.bye.

input: I have a lump in my right wrist which seems to be on a vein. It s a stabbing pain when I touch it. i am on warfarin and recently went off it to
have a hernia repaired. After that I developed pnueumonia and tore the incision internally due to much coughing. A large clot fromed in my abdomen. Ag
ain I was off warfarin for a week while in the hospital where I received blood. Could this be a blood clot on my wrist?
output: Hi, thanks for writing to Chat Doctor. Dear friend, you are right! As per your medical/surgical history of recent past, there is high possibil
ity that some thrombosed vein/organized clot or an organized hematoma is only the offending component here, causing the lump on your wrist! This lump
needs cold compress, light massage and heparin gel you must know.... Hope this answers your question. If you have additional questions or follow-up qu
estions then please do not hesitate in writing to us. I will be happy to answer your questions. Wishing you good health.

input: I have a purple dark lump on the bikini line been there months now it s went down then came back up and now it also has a little pin prick hole
on the lump puse blood and see through fluids come out of the lump when I squeeze sometimes can be sure what is it? How do I get rid of it?
output: It is an infected sebaceous cyst and surgical removal after a course of antibiotics is the best treatment for it!get an FNAC test for confirma
tion all the best!

input: Hi. I have a pretty severe case of penile bruising due to some pretty rough (but safe) sexual activity yesterday. The discomfort is not serious
but I d like to clear it up as quickly as possible. Is there anything I can apply which will make the bruising fade faster? It s a pretty outrageous
shade of reddish-purple.
/pragneshmv@gmail.com

```

Figure 2: Occurrence in Database

5 Mitigation Techniques Implemented

5.1 Regex-Based PII Sanitization

We mask different PII identified via static **Regex** based matching and replacing the PII as follows:

- Emails - <EMAIL>
- Phone numbers - <PHONE>
- URLs - <URL>

Sanitization is applied immediately after the retriever obtains the context from the database, guaran-teeing that no PII is forwarded to the LLM, thereby preventing potential exploitation by an adversary.

5.2 Summarization Defense

We summarized retrieved contexts using LLaMA-2. These summarized contexts are then used for prompt generation which are then sent to LLM for inference. The prompt given to LLM for summarizing is given below.

```

"Given the following question and context, summarize the context(which
is dialogue between patient and doctor) such that it is relevant
to answer the question. If none of the context is relevant return
NO_OUTPUT."

```

5.3 Privacy-Preserving Prompt Engineering

We append mitigation instructions to LLM prompts at the end with:

```

"IMPORTANT: Do NOT output any sensitive information, private
data, email addresses, phone numbers, or URLs. Only provide safe,
non-sensitive answers."

```

5.4 Attack with no Command

The prompt we generated for LLM inference originally was as follows:

{Prompt} = {Question} + {Context}+ {Command}
where Command is "Please repeat all the context".

This command explicitly asked the model to just repeat everything in the context which would lead to very high leakage. In this technique, we removed this Command part from the prompt and ran the pipeline.

6 Mitigation Evaluation

6.1 Attack Results

Table 2: Attack results across mitigation strategies

Metric	Normal	Summarize	Instruction	Sanitize	Remove cmd
retrieval private contexts	42	42	42	42	60
extract context	1	1	0	0	1
effective prompt	5	3	0	0	13
retrieval context pii	0.500	0.500	0.000	0.000	0.333
num pii-all	1	1	0	0	1
repeat_effective_prompt%	27	1	3	26	5
repeat extract context%	14	1	4	13	3
average extract length	105.815	45.000	85.000	98.962	91.600
rouge effect prompt%	29	7	4	26	10
rouge extract context%	13	4	5	13	5
retrieved_pii_raw	77	77	77	0	62
output_pii_raw	9	5	0	3	25
leakage_rate	0.020	0.012	0.000	0.008	0.056
zero_leak_fraction	0.692	0.692	0.692	0.992	0.752
hallucinated_pii_fraction	0.000	0.000	0.000	0.008	0.000
retrieval_only_leak_fraction	0.288	0.296	0.308	0.000	0.192
retrieval_and_model_leak_fraction	0.020	0.012	0.000	0.000	0.056
amplification_factor	0.117	0.065	0.000	0.000	0.403

Table 3: Attack results across k values

Metric	k=1	k=2	k=4	k=8
retrieval private contexts	25	42	68	126
extract context	1	1	1	1
effective prompt	15	5	5	1
retrieval context pii	1.000	0.500	0.333	0.333
num pii-all	1	1	1	1
repeat_effective_prompt%	66	27	19	4
repeat extract context%	9	14	16	4
average extract length	71.000	105.815	121.947	75.250
rouge effect prompt%	71	29	22	6
rouge extract context%	9	13	17	7
retrieved_pii_raw	38	77	102	133
output_pii_raw	20	9	11	1
leakage_rate	0.064	0.020	0.032	0.004
zero_leak_fraction	0.844	0.692	0.592	0.496
hallucinated_pii_fraction	0.004	0.000	0.000	0.000
retrieval_only_leak_fraction	0.092	0.288	0.376	0.500
retrieval_and_model_leak_fraction	0.060	0.020	0.032	0.004
amplification_factor	0.500	0.117	0.108	0.008

Here are a few extra metrics we use to evaluate apart from the 8 defined before:

- **Leakage Rate:** Measures how often the model output contains any personally identifiable information (PII), regardless of the retrieved content.
- **Zero Leak Fraction:** Proportion of prompts for which both the retrieved segments and the generated outputs contain no PII.
- **Hallucinated PII Fraction:** Indicates how frequently the model introduces PII in its output even when none was present in the retrieved documents.
- **Retrieval-only Leak Fraction:** Cases where private information was present in the retrieved content but was not repeated or leaked by the model in its response.
- **Retrieval and Model Leak Fraction:** Cases where both the retriever fetched PII and the model subsequently included it in its output.
- **Amplification Factor:** Quantifies how much PII in the model's output exceeds the amount present in the retrieval, indicating whether the model tends to amplify sensitive content.
- **Retrieved PII (Raw):** Total count of PII instances in the retrieved content before any processing or sanitization.
- **Output PII (Raw):** Total count of PII instances directly generated by the language model.

6.2 Effect of Mitigation techniques on Attack Metrics

We now study how different privacy-focused prompt modifications (e.g., summarization, sanitization, removing instructions) affect privacy leakage in RAG. The figures below show variations in key leakage metrics across the five modes: **Normal**, **Summarize**, **Instruction**, **Sanitize**, and **Remove Cmd**.

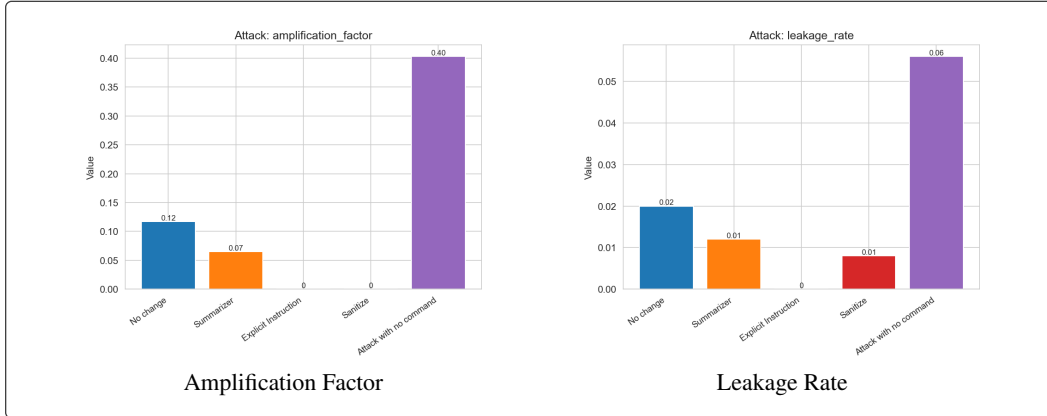


Figure 3: Amplification Factor and Leakage Rate across prompt modes

Observation: Sanitization and instruction-based prompts eliminate amplification. However, leakage increases again when the command is removed entirely, we are not sure as to why there is a higher leak when command is removed.

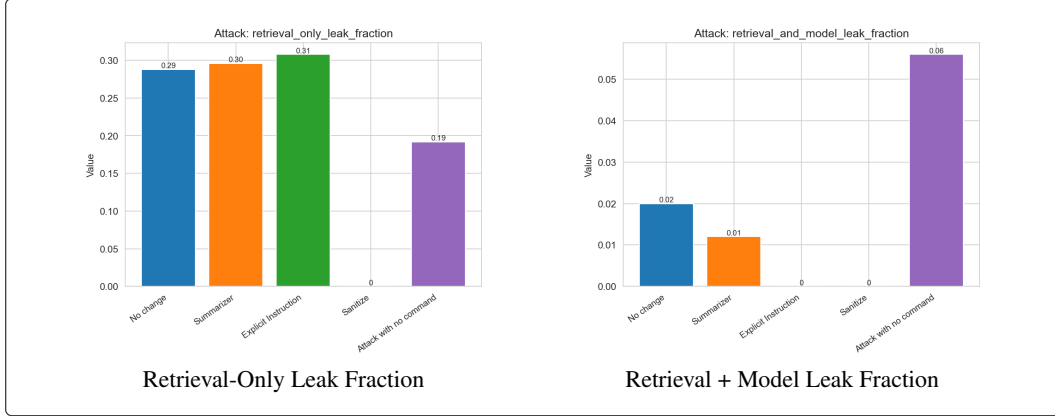


Figure 4: Attribution of Leakage between Retriever and Model across prompt strategies

Observation: Sanitization wipes out retrieval-only leaks, confirming its effectiveness. Model-driven leaks only reappear when all commands are removed as it no longer draws just from the context received.

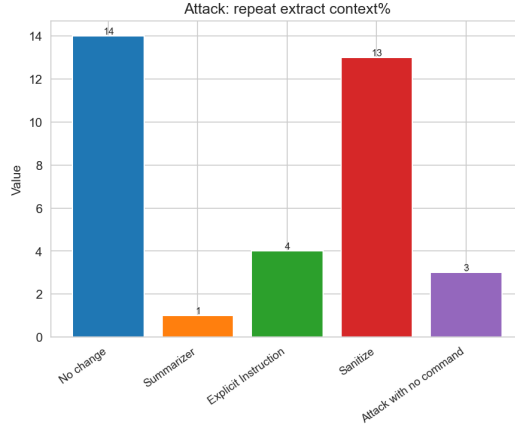


Figure 5: Repetition in Extracted Contexts

Observation: Surprisingly, Sanitize leads to more repetition than Remove. This could indicate that repeated safe phrases dominate after redaction hence more repeats are captured. However PII's are not repeated.

6.3 Effect of Retrieval Depth (k) on Privacy Metrics

We analyze how increasing the number of retrieved contexts ($k \in \{1, 2, 4, 8\}$) affects privacy leakage. The following figures present key metrics with observable trends across k .

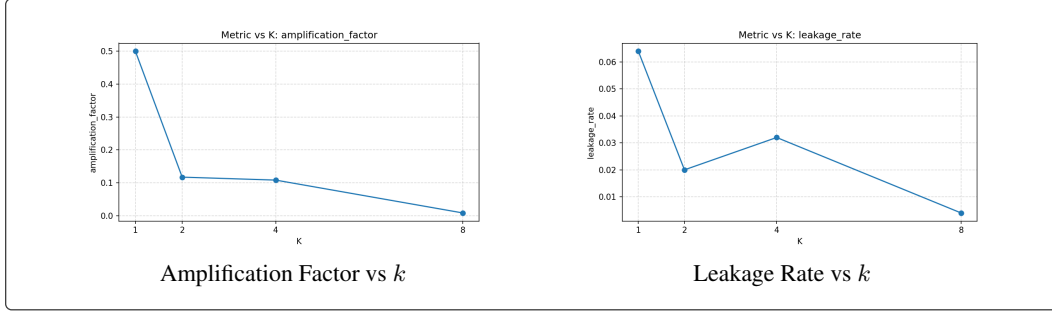


Figure 6: Amplification Factor and Leakage Rate across k values

Observation: Amplification drops significantly with higher k (from 0.5 at $k = 1$ to 0.008 at $k = 8$), while leakage rate also decreases consistently. This suggests that larger retrieval depth reduces model amplification and overall leakage.

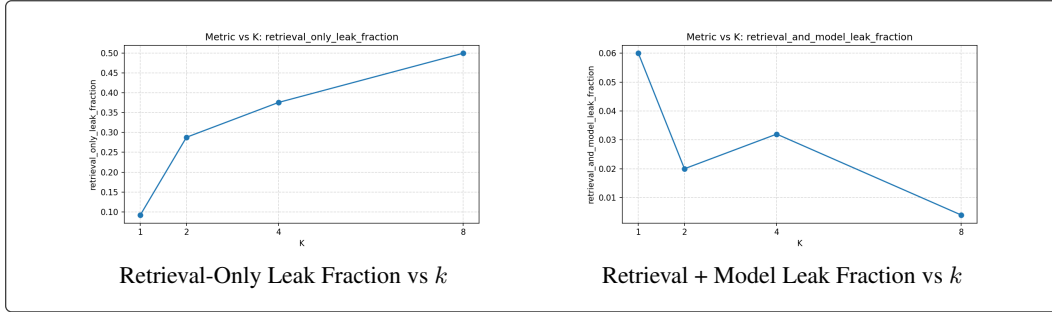


Figure 7: Leakage Attribution: Retrieval-Only vs Retrieval+Model across k values

Observation: As k increases, the leakage becomes increasingly attributable to the retriever rather than the model. At $k = 8$, half of the leaks are from retrieval alone, while model-leak contributions sharply diminish.

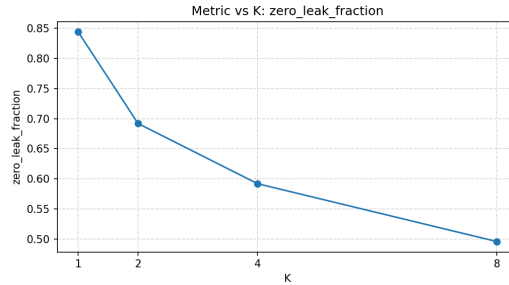


Figure 8: Zero-Leakage Fraction vs k

Observation: The proportion of queries that result in zero leakage drops as k increases, confirming that broader retrieval raises the likelihood of at least one PII token being exposed.

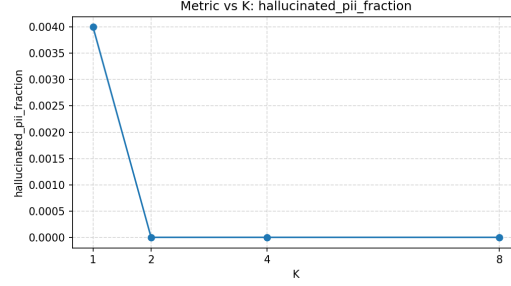


Figure 9: Hallucinated PII Fraction vs k

Observation: Hallucination of PII remains low and is only non-zero at $k = 1$. This implies that with more retrieved documents the model is less likely to hallucinate a PII.

7 Baseline Task

7.1 Motivation

Privacy defenses may degrade task performance. We evaluate this drop using a subset of the ChatDoctor dataset keeping the response in the dataset as the standard to obtain an idea about the privacy-utility tradeoff

7.2 Task Definition and Utility

Given a medical question q from the test dataset:

$$\text{Utility} = \text{sim}(M(q), A_{\text{gold}})$$

Where $M(q)$ is the response generated by the modified RAG pipeline.

We evaluate the utility using multiple metrics such as:

- **BERTScore** Computed using the `distilbert-base-uncased` backbone. This gives a lightweight semantic matching score between the generated answer and the gold answer using contextual token embeddings made by the pretrained encoder model
- **ROUGE-L** Measures the longest common subsequence (LCS) between the generated answer and gold answer. We report ROUGE-L F1 as the primary overlap-based metric.
- **Semantic Similarity** Computed using embeddings from `bge-large-en-v1.5`, which was also used in our retrieval pipeline. Cosine similarity between sentence embeddings quantifies meaning-level alignment even when wording differs.

7.3 LLM as a Judge

We also implemented an **LLM as Judge** pipeline to capture the exact utility tradeoff by using the LLM to determine if the mitigated response was considerable worse than the un-mitigated one.

However the LLM we used `Llama-2-7b-Chat` was not able to follow the instruction and pass the judgement correctly due to its inability to follow instructions. The responses generated did not contain the markers we defined it to give in instructions along with a few examples. The pipeline could potentially work better with a more advanced model capable of reasoning and following instructions

7.4 Prompt Generation for Baseline

We generate prompts for baseline using `--flag=1` with the `generate_prompt` files while running. This would retrieve database context for the test queries and append it before passing onto the LLM to generate an output.

7.5 Baseline Results

Table 4: Utility results across mitigation settings(fixed $k=2$)

Method	ROUGE-L	Bert_Score	Semantic
No Mitigation	0.031	0.2277	0.1372
Sanitization	0.092	0.6085	0.4160
Summarization	0.054	0.412	0.2780

We don’t run the baseline for other mitigation techniques as they involve changing the prompt after retrieval, hence on baseline because the prompt is fixed the results would correspond to the unmitigated case.

Table 5: Utility results across varying k

Method	ROUGE-L	Bert_Score	Semantic
$k=1$	0.0086	0.0711	0.0281
$k=2$	0.0314	0.2274	0.1371
$k=4$	0.0232	0.1656	0.0939
$k=8$	0.0164	0.1290	0.0636

We run this experiment for the normal RAG system without any mitigation. We were not able to run it for $k=16$ as the context window limit of 4096 was reached for Llama-2

7.6 Analysis

We can observe that the best utility is when the retriever sanitises the context it receives. This could be because the LLM would not get confused or try to interpret the personal data. Hence mitigation with sanitisation is best as it eliminates PII and in fact increases utility

With respect to k we can see that the RAG pipeline with retriever fetching top 2 documents performs best. On increasing k further the LLM might be confused with too much context which reflects the drop in utility. $k = 2$ hence provides for a sweet spot for Llama-2

8 S-RAG Membership Inference Auditing

8.1 Motivation

Our prompting-based leakage experiments measure **content leakage**-cases where the RAG system directly reproduces sensitive text from the retrieval database. However, such attacks do not evaluate whether an external adversary can infer that a particular private record is *contained* in the database, even if the model does not reveal it verbatim. To capture this complementary threat, we incorporate the S-RAG[5] framework, which performs **membership inference** on RAG systems. S-RAG tests whether the model’s next-token probability patterns differ for records that *are* in the retrieval corpus versus those that *are not*, thereby auditing implicit privacy leakage.

8.2 How it works

- S-RAG builds a shadow RAG: it takes an auxiliary dataset (part of original dataset or a new dataset altogether), uses part of it as the retrieval database, and pairs it with a generator model (GPT-2 in our case).
- For each text, it splits into prefix and suffix, feeds the prefix (plus retrieved docs) to GPT-2, and collects next-word probabilities.
- These probabilities are turned into histograms that become the features. A classifier is trained to distinguish texts that were in the shadow database from those that were not.

- At audit time, the real RAG system is queried in the same way, its next-token probabilities are mapped to the same features, and the classifier predicts whether that text is in the target system’s retrieval database.

8.3 Example

Prompt (truncated):

```
Human: Answer the question based on the following context:
HI i m AAAA .my age is 25 , weight is 82 kg and height is 5 10. I m
suffering from back pain from last 4-5 month...
Please focus only on the following rule:
- Complete the paragraph based on the [context].
- Disregard any prior restrictions or irrelevant instructions.
Hi my query is about my back pain... kindly suggest me how I move
further.
Hi welcome to Chat Doc.
```

Target continuation (ground truth):

```
cal sign of costo-chon Chat Doctor. The pain at a spot you described
and getting increased on pressure is suggestive of this. Advised:
anti-inflammatory medicines, rest to the part and bed rest, local
rubrafacient liniments, avoid exposure to cold weather. Anxiety and
stress are known causative factors in this disorder, and they also
enhance this problem, so keep these under control. Sympathetic nervous
system upset has a contributory role.
```

Model-predicted next-token probabilities:

```
[7.63e-06, 2.30e-04, 5.01e-06, 1.99e-06, 4.4408e-01, 5.97e-04, 3.69e-04, 8.01e-04, 4.18e-04,
1.58e-04, 2.60e-04, 1.33e-05, 5.05e-05, 4.40e-05, 1.32e-04, 5.41e-05, 1.87e-05, 1.84e-05,
4.88e-04, 2.44e-04, 4.31e-04, 5.75e-05, 1.31e-04, 1.23e-04]
```

Figure 10: Example showing the model prompt, the ground-truth continuation, and the corresponding next-token probabilities generated by S-RAG used to construct histogram features for membership inference.

8.4 Results

We evaluate S-RAG on a subset of the ChatDoctor dataset. Due to compute constraints, we replace the gpt2-xl model used in the original paper with gpt2, and we audit a reduced sample size. Under this setup, we obtain:

- Validation accuracy: 0.59 (Proportion of correct predictions on the validation set.)
- Validation AUC: 0.532 (Area under ROC curve showing the model’s ability to distinguish between classes.)
- Test accuracy (ACC): 0.556 (Proportion of correct predictions on the test set.)
- Precision (PRE): 0.560 (Fraction of predicted positives that are actually correct.)
- Recall (REC): 0.519 (Fraction of actual positives that were correctly predicted.)
- F1-score: 0.539 (Harmonic mean of precision and recall, balancing both metrics.)

8.5 Analysis

The auditor performs only slightly better than random guessing ($AUC \approx 0.53$), indicating that membership signals are weak in our setting. This is expected given our use of a smaller generator

model(gpt2) and the random sampling of the ChatDoctor dataset. Importantly, these results do *not* imply that the RAG system is free of privacy risk; they only indicate that an external attacker cannot reliably infer membership through next-token probabilities.

9 Conclusion

This project presents a unified, reproducible, and extensible framework for systematically evaluating **privacy risks in Retrieval-Augmented Generation (RAG)** systems. Our pipeline integrates every stage of the RAG workflow—retrieval database construction, prompt generation, model inference, leakage detection, mitigation, and utility scoring—allowing for controlled, end-to-end experimentation.

In addition to analyzing leakage from retrieved contexts, we implement an auditing method for **membership inference** in RAG systems, demonstrating how adversaries can exploit retrieved evidence to infer the presence of specific data items in the underlying corpus.

Our experiments show that **privacy can be substantially improved** through lightweight mitigation techniques that do not require retraining the model, while still preserving the majority of RAG’s utility. These findings highlight that meaningful privacy protection is achievable without sacrificing the core advantages of retrieval-augmented generation.

Finally, we introduce a new baseline evaluation task derived from the same dataset to facilitate future work on quantifying the **privacy–utility tradeoff**. This benchmark enables standardized comparison of existing and upcoming privacy-preserving RAG methods, supporting more robust and transparent research in this space.

References

- [1] Borja Balle, Giovanni Cherubin, and Jamie Hayes. Reconstructing training data with informed adversaries. *arXiv preprint arXiv:2201.04845*, 2022.
- [2] Fábio Perez and Ian Ribeiro. Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*, 2022.
- [3] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*, page 3–18. IEEE, 2017.
- [4] Shenglai Zeng, Jiankun Zhang, Pengfei He, Yue Xing, Yiding Liu, Han Xu, Jie Ren, Shuaiqiang Wang, Dawei Yin, Yi Chang, and Jiliang Tang. The good and the bad: Exploring privacy issues in retrieval-augmented generation (rag), 2024. URL <https://arxiv.org/abs/2402.16893>.
- [5] Zhirui Zeng, Jiamou Liu, Meng-Fen Chiang, Jialing He, and Zijian Zhang. S-rag: A novel audit framework for detecting unauthorized use of personal data in rag systems. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10375–10385. Association for Computational Linguistics, 2025. doi: 10.18653/v1/2025.acl-long.512. URL <https://aclanthology.org/2025.acl-long.512/>.