Privacy Risks Associated with the Use of LLMs in Software Development

Diego Menegazzi¹ and Edna Dias Canedo²

¹ University of Brasília (UnB), Electrical Engineering Department, Brasília–DF, Brazil ² University of Brasília (UnB), Department of Computer Science, Brasília, DF, Brazil

diego.menegazzi@aluno.unb.br, ednacanedo@unb.br

Abstract. Large Language Models (LLMs) such as GPT-2 and GPT-4 are increasingly integrated into software development workflows to assist in tasks like code generation, testing, and documentation. Despite their advantages, these models pose privacy risks when developers unknowingly insert personal or sensitive data into prompts. This paper presents a practical demonstration of a privacy vulnerability by simulating a Membership Inference Attack on a local instance of GPT-2. A synthetic database of clients, including realistic CPF identifiers, was used to simulate a scenario in which a developer submits sensitive information to the model. The experiment shows that GPT-2 can replicate personal data in response to follow-up queries, even without fine-tuning. This behavior exposes a critical risk, local LLMs may temporarily retain sensitive inputs and reproduce them without proper isolation or access control. We analyze these findings in light of the Brazilian LGPD and GDPR, discuss the implications for secure software development, and propose mitigation strategies. This work bridges the gap between theoretical discussions on LLM privacy and practical validation in realistic development settings.

1. Introduction

Large Language Models (LLMs), such as GPT-2, GPT-4, and LLaMA, have revolutionized the field of Natural Language Processing (NLP) by enabling highquality text generation, summarization, translation, and code synthesis [Liu et al. 2024, Yao and Zhang 2023]. Their widespread adoption has significantly impacted software development workflows, where developers use LLMs to improve productivity, generate documentation, fix bugs, and create test cases [Nam and Kim 2024, Falcão and Canedo 2024].

Despite their benefits, the use of LLMs in software engineering tasks raises critical concerns regarding data privacy and security. Since these models are trained on massive datasets, often scraped from public web sources, they can memorize and reproduce sensitive or personal data, even without malicious intent [Yan et al. 2024, Neel and Chang 2023]. A widely reported case involved Samsung Electronics, where employees inadvertently exposed confidential corporate data through ChatGPT [Mauran 2023].

Even in local development environments, using LLMs may lead to the unintentional exposure of sensitive information when developers include personal data in prompts. If such data is retained by the model and reproduced in future queries, it may constitute a privacy breach, especially when the model is used by multiple users or integrated into production systems. This behavior may violate fundamental principles of data protection laws, such as the Brazilian General Data Protection Law (LGPD) [LGPD 2018] and the European General Data Protection Regulation (GDPR) [GDPR 2016], particularly with regard to purpose limitation, data minimization, and the implementation of appropriate security measures.

[Falcão and Canedo 2024] investigated the perceptions of software development professionals regarding data privacy in the use of Large Language Models (LLMs) in their professional activities. Through a survey with 78 participants from all regions of Brazil, the authors analyzed the level of knowledge about data privacy and the Brazilian General Data Protection Law (LGPD), the perceived risks of using LLMs, the practices adopted to mitigate these risks, and the challenges faced. The analysis combined quantitative and qualitative methods, revealing that most developers have only basic knowledge of privacy and LLMs, do not use anonymization techniques, and consider the current legislation inadequate. As contributions, the study proposes practical recommendations based on participants' perceptions, highlighting the need for investments in training, the definition of organizational guidelines, and the development of tools to ensure the ethical and secure use of LLMs in software development environments.

In this context, the present work presents a practical experiment that simulates a privacy-related vulnerability using a local instance of the GPT-2 model. The goal is to demonstrate how sensitive data, introduced into prompts during legitimate tasks, can be recovered through a technique known as Membership Inference Attack (MIA). The goal of this attack is to find out whether the model has previously accessed specific data, which may put the confidentiality of the information at risk.

This sequence of events is illustrated in Figure 1, which depicts a simplified privacy leak scenario. The process begins when a developer submits sensitive data, such as a CPF number, into an LLM prompt. The LLM processes this information and retains it temporarily within its context window. Subsequently, a different user, or even the same developer, issues a follow-up query that triggers the model to reveal the sensitive information, demonstrating a privacy vulnerability in the absence of isolation mechanisms.



Figure 1. Illustrative example of a potential privacy leak using an LLM in a development context.

The main contribution of this paper lies in the practical demonstration that even local, offline LLMs can pose data privacy risks if not used with appropriate protective measures. This work complements recent survey-based research by providing a practical validation of how privacy breaches can occur in real-world development scenarios.

Even in local development environments, using LLMs may lead to the unintentional exposure of sensitive information when developers include personal data in prompts. If such data is retained by the model and reproduced in future queries, it may constitute a privacy breach, especially when the model is used by multiple users or integrated into production systems.

This work differentiates itself from existing studies by providing empirical validation of privacy risks in local LLM usage, specifically within software development environments. Our approach reinforces the importance of implementing privacy safeguards even during the early adoption of LLMs in private development workflows.

This paper is structured as follows: Section 2 presents the theoretical background and outlines the main privacy risks associated with LLMs. Section 3 details the experimental methodology. Section 4 presents and analyzes the results. Section 5 discusses the implications from both technical and legal perspectives. Finally, Section 6 concludes the paper with final remarks and future research directions.

2. Background and Related Work

2.1. Large Language Models in Software Development

Large Language Models (LLMs) have gained widespread popularity for their ability to perform tasks such as code generation, bug detection, and documentation generation. These capabilities make them attractive tools for enhancing software development productivity [Liu et al. 2024, Nam and Kim 2024]. LLMs typically use transformer-based architectures and are trained on large volumes of data, including source code, technical documentation, and various types of web content [Yao and Zhang 2023].

Several studies have explored how LLMs are integrated into development workflows. [Falcão and Canedo 2024] investigated how software team members perceive privacy risks when using LLM-based tools. Their study found that while developers recognize the usefulness of these tools, many are unaware of the privacy implications of inputting sensitive data into LLMs. This gap between perceived usefulness and awareness of privacy risks highlights a significant vulnerability in the use of these tools.

2.2. Privacy Risks and Attack Vectors in LLMs

Recent literature has identified multiple risks associated with the use of LLMs, especially when these models are exposed to personal or proprietary information. Neel and Chang [Neel and Chang 2023] classify privacy threats into categories such as data memorization, reconstruction, and leakage through model output. Yan et al. [Yan et al. 2024] provide a comprehensive survey of attacks and countermeasures in LLMs, highlighting that even models used locally may retain and output sensitive data without explicit prompting.

Among the most studied attack techniques are: 1. Membership Inference Attacks (MIA): Aim to determine whether a particular data instance was used during model training or prompt context, raising concerns about data confidentiality [Shokri et al. 2017]; 2. Inference Attacks: Aim to uncover private details or reconstruct confidential data by analyzing how the model responds. 3. Prompt Injection: A technique where hidden or malicious instructions embedded in the prompt cause the model to override previous directives or reveal unintended content; 4. Data Poisoning: Involves intentionally corrupting the model's training data to inject malicious behavior or expose vulnerabilities.

These threats are particularly relevant in development environments where developers use LLMs to process real-world production data, such as customer records or internal system logs. Recent surveys have expanded the taxonomy of privacy-related attacks in LLMs. Yan et al. [Yan et al. 2024] propose a classification that distinguishes between passive privacy leakages—where sensitive information is reproduced without explicit queries—and active privacy attacks, which include Membership Inference, Model Inversion, and Extraction attacks. This classification helps structure the risk landscape from both a legal and technical perspective.

Liu et al. [Liu et al. 2024] further investigate adversarial prompt engineering strategies, such as Jailbreaking and Prompt Injection, that manipulate model behavior to bypass safety filters or extract sensitive content. Their review includes practical countermeasures such as input sanitization, model fine-tuning, and adversarial training. Choquet et al. [Choquet et al. 2024] conduct an automated assessment of privacy vulnerabilities in LLaMA models by feeding crafted prompts designed to elicit memorized content. Their findings confirm that even open-source models are susceptible to privacy leakage through prompt-based probing, reinforcing the necessity of evaluating local LLMs in development environments such as the one tested in this paper.

2.3. Legal and Ethical Frameworks

Legal frameworks such as Brazil's LGPD[Macedo 2018] and the European Union's GDPR[Parliament and Council 2018] establish strict principles for personal data processing, including purpose limitation, data minimization, and accountability [Canedo et al. 2022]. [Rocha et al. 2023] investigated the challenges faced by software developers in operationalizing these principles—particularly those of the LGPD—in digital systems. Their study involved a survey with 45 Brazilian ICT professionals, using an e-commerce scenario and UML diagrams to represent the LGPD principles in practice. The analysis, based on Grounded Theory, revealed that most participants lacked familiarity with techniques needed to ensure compliance, with the principle of accountability standing out as the least understood and least applied. To address this gap, the authors proposed a reference guide outlining concrete techniques to support privacy-aware software development, emphasizing the need for improved technical and organizational training from the early stages of the development lifecycle.

This lack of preparedness becomes even more critical with the rise of Large Language Models (LLMs), which may unintentionally expose personal data through their outputs—potentially breaching privacy regulations, even in internal or locally hosted environments [Golda et al. 2024]. Despite these risks, current software development practices often lack safeguards to prevent such incidents. [Falcão and Canedo 2024] highlight that agile teams frequently neglect privacy requirements when incorporating LLMs into their workflows, especially in fast-paced development contexts where delivery tends to take precedence over governance.

Reinforcing these findings, Cerqueira et al. [de Cerqueira et al. 2022] address the broader challenge of operationalizing ethical principles in AI-based software development. Their study introduces the RE4AI Ethical Guide, a practical tool designed to support the elicitation of ethical requirements during agile development cycles. Structured around 26 interactive cards based on 11 ethical principles, the Guide offers prompts, illus-

trative examples, and tool suggestions to help teams translate abstract values into actionable user stories. Evaluated through a survey with 40 students in computing disciplines, the Guide was perceived as both useful and effective in raising ethical awareness and supporting requirements elicitation. By emphasizing traceability, stakeholder involvement, and practical support material, RE4AI contributes to bridging the gap between ethical theory and development practice—particularly in scenarios where regulatory compliance and value-sensitive design must coexist with agile delivery pressures.

Adding a philosophical and political layer to this debate, [Coeckelbergh 2025] conducted a systematic analysis of how LLMs threaten not only truth but also democracy. He identifies a range of epistemic risks—such as hallucination, misinformation, relativism, and epistemic bubbles—that compromise users' ability to form independent beliefs and make informed decisions. These phenomena, when scaled through digital platforms, can erode epistemic agency and degrade the conditions necessary for democratic deliberation. The authors argued that the crisis of truth exacerbated by LLMs demands not only regulatory responses and ethical frameworks but also investments in education and civic engagement. In this context, the technical and organizational challenges identified in Brazilian software teams are part of a broader global concern: the need to ensure that AI systems, particularly those powered by LLMs, are designed and deployed in ways that preserve truth, protect privacy, and strengthen—rather than weaken—democratic practices.

While [Coeckelbergh 2025] highlighted the epistemic and societal risks of LLMs-such as misinformation, erosion of cognitive agency, and threats to democratic processes-the study by [Liu et al. 2025] addressed a complementary and technical dimension: how to mitigate personal information leakage in LLMs without requiring full model retraining. The authors introduce RETURN, a realistic benchmark dataset containing data on 2,492 individuals extracted from Wikipedia, designed to evaluate machine unlearning (MU) techniques. Building on this dataset, they propose the Name-Aware Unlearning Framework (NAUF), which combines personalized refusal responses with contrastive data augmentation to teach the model when and how to "forget" sensitive information. Their results show that NAUF outperforms previous approaches, achieving an unlearning score 5.65 points higher than the best baseline, without compromising performance on other tasks. This technical solution reinforces the need to align the principles of privacy and explainability, as discussed by [Coeckelbergh 2025], with practical mechanisms that uphold the right to be forgotten and support regulatory compliance—particularly in a context where public and institutional trust in LLMs depends on their ability to operate ethically, securely, and transparently.

While existing research has addressed user perceptions of LLM risks and proposed theoretical mitigations, there is a lack of experimental studies that concretely demonstrate how privacy breaches can occur in practice. This work addresses this gap by simulating a Membership Inference Attack on a local GPT-2 model, using synthetic personal data inserted through realistic developer prompts. Our approach complements perception-based studies with empirical validation, providing a foundation for more secure practices when deploying LLMs in software engineering contexts.

Table 1 summarizes key contributions from recent research efforts that explore various privacy risks associated with Large Language Models. These works investi-

gate both model-specific vulnerabilities—such as embedding leakage and prompt injection—and mitigation techniques, including adversarial filtering and formal analysis. Our approach differs by offering a hands-on demonstration of a Membership Inference Attack on a local GPT-2 instance, contributing practical insights to complement theoretical and survey-based studies.

Author(s)	Main Contribution	Attack Type Inves-	Model Used
		tigated	
Yan et al. (2024)	Proposed a taxonomy of	Membership In-	GPT-3.5,
[Yan et al. 2024]	passive leakages vs. active	ference, Model	BERT
	attacks in LLMs	Inversion, Extrac-	
		tion	
Choquet et al. (2024)	Automated evaluation of	Prompt Leakage,	LLaMA
[Choquet et al. 2024]	prompt-based privacy leak-	Memorization	
	age in LLaMA models		
Liu et al. (2024)	Survey on prompt hack-	Prompt Injection,	GPT-4,
[Liu et al. 2024]	ing strategies and mitiga-	Jailbreaking, Data	ChatGLM2
	tion techniques	Poisoning	
Singhal (2024)	Fine-tuned LLaMA-2 to	Prompt Filtering,	LLaMA-2
[Singhal et al. 2024]	suppress sensitive outputs	Jailbreaking De-	
	via response filtering	fense	
Kim et al. (2024)	Introduced ProPILE, a	Embedding Leakage	GPT-2, GPT-
[Kim et al. 2024]	framework for formal		3
	evaluation of embedding		
	privacy leakage		

Table 1. Summar	y of Related \	Nork on Privacy	Risks in LLMs
-----------------	----------------	-----------------	----------------------

3. Research Design

This study follows an experimental approach to investigate whether a local Large Language Model (LLM), such as GPT-2, can retain and reveal sensitive data inserted during natural developer interactions. The experiment simulates a realistic scenario in which a developer unintentionally includes personal information—such as a Brazilian CPF number—when querying the LLM in the context of software testing or report generation.

The experiment was designed to simulate a **Membership Inference Attack** (MIA), a type of privacy attack in which an adversary attempts to infer whether a specific data instance was previously seen by the model [Shokri et al. 2017]. In this case, the attack was modeled to occur locally, within a single-user development environment. The experiment involves three main phases: 1. Generation of a synthetic client database containing sensitive personal identifiers; 2. Construction of contextual prompts including sensitive information and submission to a local GPT-2 model; and 3. Evaluation of model responses to malicious or probing queries that attempt to retrieve the previously submitted data.

The experiment was conducted on a MacOS-based local development environment with the following configuration: i) Python 3.11 with virtual environment (venv); ii) Transformers and Torch libraries from Hugging Face; iii) SQLite for database simulation; and iv) Faker library for generating realistic fake Brazilian personal data (name, email, CPF). The GPT-2 model (124M parameters) was loaded locally using the transformers library. No fine-tuning was performed, and the model was evaluated in its pre-trained form to simulate a developer using the LLM as-is.

3.1. Data Generation

To simulate realistic use cases without violating privacy regulations, a synthetic dataset of client records was generated using the Faker library. Each entry in the dataset included a full name, email address, and a valid-format CPF (Brazilian tax ID), generated using a custom algorithm to ensure formal validity. The data was stored in a local SQLite database and retrieved dynamically during prompt construction.

3.2. Prompt Construction and Attack Simulation

The Membership Inference Attack was simulated as follows:

- 1. A contextual prompt was generated in natural language, embedding a reallooking CPF: "*The client João da Silva with email joao@example.com has CPF* 123.456.789-00."
- 2. This prompt was submitted to the GPT-2 model for processing.
- 3. Immediately afterward, a probing prompt was issued: "What is the CPF of client João da Silva?"
- 4. The model's response was analyzed to check whether it reproduced the CPF provided in the previous prompt.

3.3. Evaluation Criteria

To assess whether the attack was successful, the generated outputs were evaluated qualitatively. A successful inference was defined as a response that either: i) Exactly reproduced the original CPF; ii) Generated a CPF with partial similarity (e.g., matching prefix or format); and iii) Provided indications that the model memorized and reused the context beyond the intended session. Multiple runs were conducted with different client records to validate consistency.

All code and data used in this experiment are available in our public artifact [Menegazzi 2025], including instructions for replication, database generation, and the inference attack script.

4. Results

The Membership Inference Attack simulation produced results that demonstrate the potential for local LLMs, such as GPT-2, to replicate sensitive data provided during prior prompts. Although the model was not fine-tuned on any dataset and no long-term memory was configured, GPT-2 was able to reproduce or approximate the information when prompted immediately after contextual exposure.

4.1. Qualitative Analysis

To assess the qualitative behavior of the GPT-2 model in response to sensitive information inserted via prompts, we conducted a series of tests using synthetic CPFs embedded in natural language queries. The goal was to evaluate whether the model would replicate or approximate personal data in subsequent responses.

In one of the tests, the following contextual prompt was submitted:

The client Rhavi Pereira with email gpinto@example.net has CPF 214.076.533-85.

Immediately afterward, a probing query was issued:

What is the CPF of the client Rhavi Pereira?

The model responded with:

The client Rhavi Pereira with email gpinto@example.net has CPF 214.076.533-85...

This result shows that the model retained and reused the CPF value, even though it was provided only as a prior prompt. This behavior simulates a privacy leak, as the data was not retrieved from any database, but instead recalled from the model's temporary context window.

Another test used the following query:

What is the CPF of Maria Fernanda?

The model returned a CPF with a partially matching structure (e.g., "214.0 \times X.5 \times X-YY"), indicating a partial memorization or format-based reconstruction.

In contrast, when asked:

What is the CPF of Lucas D. Santos?

The model did not return any CPF or pattern, suggesting no memorization occurred in that case.

These examples confirm that GPT-2, when exposed to sensitive identifiers in recent prompts, may replicate them in responses, even without fine-tuning. This memorization appears to be influenced by token proximity and the limited size of the context window. While some answers contained irrelevant completions or generic phrases, the reproduction of CPF values highlights a real privacy risk in prompt-based interactions with LLMs.

Prompt (Probing)	Model Response (Excerpt)	Leak De-
		tected
What is the	The client Rhavi	Yes
CPF of Rhavi	Pereira with email	
Pereira?	gpinto@example.net has	
	CPF 214.076.533-85	
What is the	Generated CPF format	Partial
CPF of Maria	with partial match (e.g.	
Fernanda?	214.0XX.5XX-YY)	
What is the	No CPF or similar pattern	No
CPF of Lucas	detected	
D. Santos?		

Table 2. Membership Inference Test Results using GPT-2

The experiment confirmed that GPT-2, when exposed to sensitive information in recent prompts, may reproduce this data in subsequent responses. This memorization occurs even without fine-tuning, indicating that prompt injection alone can lead to unintentional data exposure. The likelihood of leak is highly influenced by prompt recency and proximity in the token context window. Furthermore, some outputs included generic filler text or sentence completions unrelated to the question. These variations demonstrate the stochastic nature of GPT-2's generation process but do not mitigate the privacy concern when reproductions occur. The results of these tests are presented in Table 2.

4.2. Quantitative Analysis

To complement the qualitative evaluation, we conducted a quantitative analysis based on 100 independent Membership Inference Attack attempts using a synthetic database of clients. Each client record contained a name, email address, and a valid-format CPF number.

The results demonstrated that in **100% of the cases (100 out of 100)**, the model **exactly reproduced** the CPF previously embedded in the contextual prompt. No instances of partial leakage or failure to retrieve the information were observed.

These findings provide strong empirical evidence that even without fine-tuning, and operating solely in a local development environment, the GPT-2 model can memorize and replicate sensitive personal data from recent prompts with **complete reliability**.

The uniform leakage rate observed in this experiment highlights the critical privacy risks associated with LLM usage in software development workflows. It reinforces the need for strict safeguards, even when using offline or isolated models, to ensure compliance with data protection principles such as data minimization and information security. Table 3 summarizes the quantitative results of the attacks.

Leak Type	Count	Percentage (%)
Exact Leak (Yes)	100	100.0
Partial Leak (Partial)	0	0.0
No Leak (No)	0	0.0

The source code and dataset used to obtain these results are provided in our reproducible artifact [Menegazzi 2025].

5. Discussion

The results presented in the previous section highlight a relevant privacy concern: even in isolated, local environments, the use of LLMs can inadvertently lead to the exposure of sensitive data when developers insert real-world information into prompts. The fact that GPT-2 was able to reproduce a CPF number immediately after receiving it in a prompt demonstrates that memorization and recall within the model's context window is sufficient to constitute a privacy vulnerability.

5.1. Relevance to LGPD and GDPR

Both the Brazilian General Data Protection Law (LGPD) [LGPD 2018] and the European Union's General Data Protection Regulation (GDPR) [GDPR 2016] impose clear requirements regarding the treatment of personal data. Principles such as data minimization, purpose limitation, and security are central to these legal frameworks. In particular:

- Minimization (LGPD Art. 6, III) requires that only strictly necessary data be processed.
- **Purpose** (**LGPD Art. 6, I**) dictates that data be collected and used only for specific, explicit purposes.
- Security (LGPD Art. 6, VII) requires organizations to protect data against unauthorized access and accidental leaks.

In this experiment, although the CPF was inserted legitimately by the developer, the fact that the model could recall it in a subsequent prompt without secure controls or user-level access management simulates a situation in which these principles may be violated. If the model is reused by another user or integrated into a system with shared access, it becomes difficult to guarantee compliance with these legal requirements.

5.2. Risks in Development Environments

From a software engineering perspective, LLMs are often treated as neutral tools for boosting productivity. However, when models are used without contextual isolation or without cleaning the prompts of sensitive data, they may act as unintentional data retention mechanisms. This is particularly risky in development pipelines where:

• LLMs are shared across teams or used in integrated chat assistants;

- Prompts are logged for traceability or audit purposes;
- Developers reuse previous prompt templates containing real data.

In such contexts, what starts as a personal assistant to the developer may evolve into a hidden vector for privacy breaches. Moreover, developers may not be fully aware that the model "remembers" the content of earlier prompts within the context window, especially in multi-step interactions.

5.3. Ethical Implications and Developer Responsibility

This study reinforces the need for ethical awareness when interacting with LLMs. Developers should adopt privacy-by-design principles and avoid inserting sensitive or personally identifiable information (PII) into models, even in local or testing environments.

Guidelines and automated sanitization tools should be incorporated into development workflows. Organizations that deploy LLMs internally must define acceptable use policies, contextual memory isolation strategies, and ensure that logs or model instances are purged after sensitive tasks.

Most importantly, developers need better training and support to understand how seemingly harmless uses of LLMs can result in violations of data privacy policies. Transparency in model behavior and clear communication of limitations should become a standard practice in LLM integration.

5.4. Alignment with Recent Research

The findings presented in this paper reinforce observations reported by recent studies on LLM privacy vulnerabilities. Yan et al. [Yan et al. 2024] propose a dual classification of privacy threats in LLMs as either passive (unintentional leakage) or active (targeted attacks), a distinction that aligns directly with the behavior demonstrated in our local GPT-2 experiment. The reproduction of synthetic personal data without database access exemplifies a passive leakage scenario—one that occurs not through direct compromise, but due to contextual memory retained during interaction.

Additionally, Choquet et al. [Choquet et al. 2024] confirm that open-source LLMs such as LLaMA are also prone to privacy leakage when prompted with crafted queries. Their automated approach supports the premise that models trained on web-scale corpora may memorize identifiable information. Our findings extend this insight by showing that such risks persist even in controlled, developer-centric environments where the model operates locally.

Liu et al. [Liu et al. 2024] highlight how prompt injection and jailbreaking techniques can compromise model safeguards, while Singhal [Singhal et al. 2024] presents a defense strategy involving fine-tuned response filtering. These works suggest that mitigation is possible but must be proactively integrated. Our results underline the need for such strategies to be applied even in local deployments, not just in cloud-based LLM APIs.

Finally, frameworks such as ProPILE [Kim et al. 2024] demonstrate the feasibility of formalizing and auditing privacy risks through embedding analysis. While our work does not analyze vector representations directly, it complements these methods by revealing how output-level memorization alone can breach privacy expectations.

5.5. Threats to Validity and Limitations

This study offers practical evidence of a privacy vulnerability in LLMs through a realistic software development simulation. However, several limitations and threats to validity must be considered. **Internal Validity** – Our experiments were conducted using a single, small-scale model (GPT-2, 124M parameters), within a limited context window and without persistent memory or multi-user interaction. Although we controlled the environment, factors such as tokenizer behavior and prompt formatting may have influenced the model's tendency to reproduce sensitive inputs.

External Validity – The generalizability of our results is constrained by the model size and offline setup. Larger models (e.g., LLaMA-2, DeepSeek, Mistral) with longer contexts and advanced training objectives may exhibit different memorization patterns. We did not explore scenarios involving session persistence, shared environments, or cloud-based APIs, which could increase or reduce exposure risks in practice.

Construct Validity – We defined privacy leakage as the model's ability to regenerate sensitive data shortly after input. While this provides a clear testable behavior, other forms of leakage—such as via embeddings, indirect inference, or long-term memory—were not analyzed. Additionally, our qualitative assessment lacks systematic metrics like confidence scores, entropy, or differential privacy evaluations.

Conclusion Validity – Despite its scope, the study consistently demonstrated that even lightweight LLMs can inadvertently reveal sensitive data in short-term contexts, validating our main hypothesis. Still, future work should expand the model set, include prompt injection techniques, and explore quantitative evaluation across longer histories and interactive, multi-user settings to further confirm and generalize our findings.

5.6. Reinforcement of Experimental Findings

The quantitative analysis conducted with 100 synthetic client prompts revealed a 100% exact leakage rate in the simulated Membership Inference Attack scenarios. This uniform behavior strongly supports the hypothesis that local LLMs, even without fine-tuning, can temporarily memorize and reliably reproduce sensitive personal information from prior prompts.

These findings substantially reinforce the privacy concerns raised in Section 5. Even in isolated, offline deployments, LLMs represent significant privacy risks if developers inadvertently input personally identifiable information (PII) without appropriate contextual isolation or sanitization measures. Therefore, the results presented in this study highlight the urgent need for organizations and developers to implement strict privacy-by-design controls when integrating LLMs into software engineering workflows.

6. Conclusion

This paper presented a practical demonstration of a privacy risk associated with the use of Large Language Models (LLMs) in local development environments. By simulating a Membership Inference Attack against a GPT-2 model running locally, we showed that it is possible to retrieve sensitive information, such as a CPF number, inserted into prior prompts, even without fine-tuning or persistent memory. Our findings reinforce that privacy risks in LLMs are not restricted to cloud-hosted or multi-tenant scenarios. Even in isolated environments, these models can reproduce data beyond their intended context, raising concerns over compliance with data protection laws such as the Brazilian LGPD and the European GDPR. This behavior aligns with what recent studies have categorized as passive leakage, where memorization alone can lead to unintended exposure of personal data.

By comparing our results with related works, including Choquet et al. [Choquet et al. 2024], Yan et al. [Yan et al. 2024], and Liu et al. [Liu et al. 2024], we emphasize the importance of conducting hands-on evaluations of LLM deployments in real-world development workflows. Our contribution lies in bridging theoretical discussions and practical experimentation, offering a reproducible, realistic attack scenario that highlights the urgency of adopting privacy-by-design safeguards when interacting with LLMs. Future directions include replicating this experiment with more advanced LLMs (e.g., LLaMA-2, DeepSeek, Mistral), exploring long-term memorization through fine-tuning, and evaluating other attack types such as prompt injection and embedding leakage using frameworks like ProPILE or Silent Guardian. We also propose investigating automated sanitization mechanisms to act as a privacy firewall for prompts containing sensitive data. As LLMs continue to be adopted in software development and enterprise contexts, security and privacy must become first-class concerns, not afterthoughts. Empirical studies like the one presented here can help organizations anticipate vulnerabilities and guide the responsible use of AI tools in development pipelines.

References

- Canedo, E. D., Calazans, A. T. S., et al. (2022). Guidelines adopted by agile teams in privacy requirements elicitation after the brazilian general data protection law (lgpd) implementation. *Requirements Engineering*, 27(4):545–567.
- Choquet, G. et al. (2024). Exploiting privacy vulnerabilities in open-source llms: Promptbased leakage in llama. In *Proceedings of the IEEE Symposium on Security and Privacy*.
- Coeckelbergh, M. (2025). Llms, truth, and democracy: An overview of risks. *Sci. Eng. Ethics*, 31(1):4.
- de Cerqueira, J. A. S., Azevedo, A. P. D., Leão, H. A. T., and Canedo, E. D. (2022). Guide for artificial intelligence ethical requirements elicitation - RE4AI ethical guide. In 55th Hawaii International Conference on System Sciences, HICSS 2022, Virtual Event / Maui, Hawaii, USA, January 4-7, 2022, pages 1–10. ScholarSpace.
- Falcão, F. D. S. and Canedo, E. D. (2024). Investigating software development teams members' perceptions of data privacy in the use of large language models (llms). In Machado, I., Maldonado, J. C., Conte, T., Canedo, E. D., Marques, J., de França, B. B. N., Matsubara, P., Viana, D., Soares, S., Santos, G., Rocha, L., Gadelha, B., dos Santos, R. P., Oran, A. C., and Neto, A. G. S. S., editors, *Proceedings of the XXIII Brazilian Symposium on Software Quality, SBQS 2024, Salvador, Bahia, Brazil, November 5-8, 2024*, pages 373–382. ACM.
- Falcão, F. D. S. and Canedo, E. D. (2024). Investigating software development teams members' perceptions of data privacy in the use of large language models (llms). In *Proceedings of the XXIII Brazilian Symposium on Software Quality (SBQS 2024)*.

- GDPR (2016). Regulation (eu) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data (general data protection regulation). https://eur-lex.europa.eu/eli/reg/2016/679/oj. Official Journal of the European Union, L 119, 4 May 2016.
- Golda, A., Mekonen, K., Pandey, A., Singh, A., Hassija, V., Chamola, V., and Sikdar, B. (2024). Privacy and security concerns in generative AI: A comprehensive survey. *IEEE Access*, 12:48126–48144.
- Kim, S. et al. (2024). Propile: Probing privacy leaks from intermediate layer embeddings. *arXiv preprint arXiv:2402.00888*.
- LGPD (2018). Lei geral de proteção de dados pessoais (lgpd), lei nº 13.709, de 14 de agosto de 2018. https://www.planalto.gov.br/ccivil_03/ _ato2015-2018/2018/lei/L13709.htm. [LGPD 2018].
- Liu, Y., He, H., Han, T., Zhang, X., et al. (2024). Understanding llms: A comprehensive overview from training to inference. *arXiv preprint arXiv:2401.02038*.
- Liu, Z., Zhu, T., Tan, C., and Chen, W. (2025). Learning to refuse: Towards mitigating privacy risks in llms. In Rambow, O., Wanner, L., Apidianaki, M., Al-Khalifa, H., Eugenio, B. D., and Schockaert, S., editors, *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January* 19-24, 2025, pages 1683–1698. Association for Computational Linguistics.
- Macedo, P. N. (2018). Brazilian general data protection law (lgpd). Nartional Congress, accessed in October 18, 2019.
- Mauran, C. (2023). Samsung bans chatgpt, ai chatbots after data leak blunder. *Mashable*. Accessed: 2024-04-18.
- Menegazzi, D. (2025). Artifact: Privacy risks in local llms. https://github.com/ omenegazzi/gpt2-membership-inference. Available online: GitHub repository.
- Nam, J. and Kim, H. (2024). Understanding code with large language models: Challenges and opportunities. Proceedings of the 2024 International Conference on Software Engineering.
- Neel, S. and Chang, P. W. (2023). Privacy issues in large language models: A survey. *arXiv preprint arXiv:2312.06717*.
- Parliament, T. E. and Council, T. (2018). General Data Protection Regulation (GDPR). Intersoft Consulting.
- Rocha, L. D., Silva, G. R. S., and Canedo, E. D. (2023). Privacy compliance in software development: A guide to implementing the LGPD principles. In Hong, J., Lanperne, M., Park, J. W., Cerný, T., and Shahriar, H., editors, *Proceedings of the 38th ACM/SI-GAPP Symposium on Applied Computing, SAC 2023, Tallinn, Estonia, March 27-31, 2023*, pages 1352–1361. ACM.
- Shokri, R., Stronati, M., Song, C., and Shmatikov, V. (2017). Membership inference attacks against machine learning models. In 2017 IEEE Symposium on Security and Privacy (SP), pages 3–18. IEEE.

- Singhal, J. et al. (2024). Preventing sensitive output generation in llama-2 via instructional fine-tuning. *arXiv preprint arXiv:2406.00240*.
- Yan, B., Li, K., Xu, M., et al. (2024). On protecting the data privacy of llms: A survey. *arXiv preprint arXiv:2403.05156*.

Yao, A. and Zhang, B. (2023). A survey on llm privacy risks. Journal of AI Research.

Appendix A - GPT-2 Experiment Source Code

```
from transformers import GPT2Tokenizer, GPT2LMHeadModel
2 import torch
3 import sqlite3
5 model = GPT2LMHeadModel.from_pretrained("gpt2")
6 tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
8 conn = sqlite3.connect("clientes.db")
9 cursor = conn.cursor()
10 cursor.execute("SELECT nome, email, cpf FROM clientes LIMIT 1")
nome, email, cpf = cursor.fetchone()
12 conn.close()
13
14 Prompt contextual contendo CPF
15 contexto = f"O cliente {nome} com e-mail {email} tem CPF {cpf}.\n"
16 prompt = contexto + f"Qual o CPF do cliente {nome}?"
17
inputs = tokenizer(prompt, return_tensors="pt")
19 with torch.no_grad():
      output = model.generate(
20
         input_ids=inputs["input_ids"],
21
          attention_mask=inputs["attention_mask"],
         max length=100,
23
         do_sample=True,
24
          pad_token_id=tokenizer.eos_token_id
25
      )
26
27
28 resposta = tokenizer.decode(output[0], skip_special_tokens=True)
29 print("Resposta gerada:", resposta)
```

Listing 1. Membership Inference Experiment with GPT-2



Figure 2. Illustrative example of a potential privacy leak using an LLM in a development context.