

Whispered Tuning: Data Privacy Preservation in Fine-Tuning LLMs through Differential Privacy

Tanmay Singh¹, Harshvardhan Aditya¹, Vijay K. Madiseti², Arshdeep Bahga³

¹School of Computer Science Engineering & Technology, Bennett University, Greater Noida, India

²School of Cybersecurity and Privacy, Georgia Institute of Technology, Atlanta, USA

³Cloudemy Technology Labs, Chandigarh, India

Email: singhtanmay0915@gmail.com, harshva27@gmail.com, madiseti.vijay@gmail.com, arshdeep@cloudemy.io

How to cite this paper: Singh, T., Aditya, H., Madiseti, V.K. and Bahga, A. (2024) Whispered Tuning: Data Privacy Preservation in Fine-Tuning LLMs through Differential Privacy. *Journal of Software Engineering and Applications*, 17, 1-22.
<https://doi.org/10.4236/jsea.2024.171001>

Received: December 1, 2023

Accepted: January 20, 2024

Published: January 23, 2024

Copyright © 2024 by author(s) and Scientific Research Publishing Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

The proliferation of Large Language Models (LLMs) across various sectors underscored the urgency of addressing potential privacy breaches. Vulnerabilities, such as prompt injection attacks and other adversarial tactics, could make these models inadvertently disclose their training data. Such disclosures could compromise personal identifiable information, posing significant privacy risks. In this paper, we proposed a novel multi-faceted approach called Whispered Tuning to address privacy leaks in large language models (LLMs). We integrated a PII redaction model, differential privacy techniques, and an output filter into the LLM fine-tuning process to enhance confidentiality. Additionally, we introduced novel ideas like the Epsilon Dial for adjustable privacy budgeting for differentiated Training Phases per data handler role. Through empirical validation, including attacks on non-private models, we demonstrated the robustness of our proposed solution SecureNLP in safeguarding privacy without compromising utility. This pioneering methodology significantly fortified LLMs against privacy infringements, enabling responsible adoption across sectors.

Keywords

NLP, Differential Privacy, Adversarial Attacks, Informed Decisions

1. Introduction

The rapid development of Large Language Models (LLMs) by organizations such as Meta and OpenAI, exemplified by models like Llama 2 [1] and GPT-4, marks a significant advancement in the AI field. These models are benchmarks in per-

formance and security, but their widespread use also brings critical challenges, particularly in data privacy. This is a concern not only for models backed by large institutions but also for those developed by independent researchers or smaller groups, where safety measures might be less stringent.

LLMs, with their expansive data processing capabilities, are prone to vulnerabilities that could lead to inadvertent data exposure. This is especially troubling in sectors handling sensitive information, such as healthcare, finance, and law. The risk of revealing personal or confidential data through advanced techniques like prompt injection attacks [2] is a real and present danger.

In our paper, we explore the intrinsic vulnerabilities of LLMs in depth, focusing on the mechanisms of potential data leakage. This analysis is crucial in understanding the limitations these models face when dealing with sensitive data, a step vital for creating more secure AI systems.

The fine-tuning of LLMs with sensitive data adds to the complexity of ensuring data privacy. This process increases the risk of embedding private details into the model's parameters, making it challenging to ensure that the final model is devoid of sensitive information. While methods like deduplication of training data can enhance security, they are not foolproof.

Our proposed solution, SecureNLP, addresses these privacy concerns. It integrates a PII redaction model with differential privacy mechanisms to strike a balance between the utility of LLMs and the need to protect sensitive information. The core of our approach includes Whispered Tuning, a novel training method that enhances resistance to data extraction attempts, and the application of differential privacy to prevent the retention of identifiable data patterns.

By addressing these challenges, we aim to contribute to the development of LLMs that are not only powerful in processing and generating information but also robust in terms of security and privacy, ensuring their safer application in various critical domains.

This paper presents a multi-faceted Whispered Tuning approach integrating PII redaction, differential privacy, output filtering, and architectural improvements to enhance privacy preservation in LLMs. The key research contributions and novel aspects of the paper are as follows: Proposes a 4-step approach called “Whispered Tuning” to address training data privacy leaks in large language models (LLMs)—Uses a PII redaction model in the first step to redact personally identifiable information (PII) from the dataset—Fine-tunes the models on the redacted dataset in the second step—Introduces differential privacy techniques during fine-tuning using the Opacus library—Implements a “Self Reflection Filter” to check model outputs for any remaining PII and replace with fake data if found—Compares models fine-tuned with and without PII redaction (“Clear-View Model” vs “SecureNLP”)—Demonstrates vulnerability of ClearView Models to “canary insertion attacks” which extract private data—Proposes new ideas over baseline like the Epsilon Dial and differentiated Training Phases by Role.

In this paper, we commence with a detailed Statement of the Problem, setting

the stage for our research by outlining the core issue and providing necessary background information. Following this, we delve into the Whispered Tuning: Architecture section, which is subdivided into critical components: Firstly, the Dataset Generation & Analysis section, with a special focus on Exploratory Data Analysis, lays the groundwork for our study by detailing the dataset preparation and its comprehensive examination. This is followed by a discussion on PII Classification methods used, and the process of Fine-Tuning on a Redacted Dataset, highlighting our approaches to enhancing data privacy. We also explore the novel concept of Fine-Tuning with a Private Optimizer and introduce a unique Filter for Output—Self Reflexion Prompt to further secure model outputs. The Main Results section presents our pivotal findings, including the Fine-Tuning of different models on a Dataset with PII to create ClearView Models, primarily to establish a benchmark, an analysis of Different Types of Attack on ClearView Models, and a Comparison between ClearView Models and SecureNLP Models: Private LLMs. Finally, the paper concludes with the Conclusions and Suggestions section, where we synthesize our findings and offer forward-looking recommendations and implications derived from our research.

2. Related Work

While their ability to process and generate information has transformed various sectors, it has also introduced significant privacy risks, especially concerning the potential leakage of Personal Identifiable Information (PII). The literature addressing these concerns is rich and evolving, with several noteworthy contributions that lay the groundwork for our proposed solution, SecureNLP.

Weiyang Shi *et al.*'s work "Just Fine-tune Twice: Selective Differential Privacy for Large Language Models" [3] serves as a foundational pillar for our research. Their approach emphasizes the need for differential privacy in fine-tuning LLMs, a principle that we have expanded upon in SecureNLP. By selectively applying differential privacy, they offer a trade-off between model utility and privacy, which we aim to optimize further in our work.

Klymenko *et al.* [4] present a comprehensive review of differential privacy applications in natural language processing, underscoring the critical balance between data utility and privacy. This review sets the stage for the necessity of privacy-preserving techniques in LLMs, reinforcing the relevance of our approach.

Behnia *et al.* [5] investigated fine-tuning methods for Large Language Models (LLMs) using differential privacy, a study that closely mirrors our approach. Their research focuses on the intricate balance between preserving privacy and maintaining the model's performance. SecureNLP addresses this challenge by introducing innovative architectural improvements, directly tackling the trade-off between privacy and performance.

The technical report by Li *et al.* [6], "Privacy-Preserving Prompt Tuning for Large Language Model Services," presents a targeted approach to safeguard LLMs during the prompt tuning phase. While their focus is narrower, it complements

our broader goal by demonstrating the effectiveness of privacy-preserving measures in specific stages of model training.

Informed by this extensive literature and the challenges identified, our research contributes novel solutions and perspectives to the field, aiming to address these gaps in privacy and security within LLMs.

Our work significantly extends the foundational research of Shi *et al.* [3], introducing several pivotal enhancements. We have developed a state-of-the-art PII classification model, specifically for identifying and classifying Personal Identifiable Information (PII) in the “pii-masking-65k” dataset. In addition to this, our separate adversarial attack investigations have uncovered vulnerabilities in a variety of Large Language Models (LLMs). A major breakthrough in our research is the introduction of the “Self-Reflection Filter,” where model outputs are re-analyzed by the PII classification model to ensure the effective removal of any inadvertent PII leakages. We have also suggested the “Epsilon Dial,” an innovative concept for adjustable privacy budgeting, and have redefined Training Phases by Role, optimizing the training process based on the data handler’s role. Importantly, our enhanced architectural approach is tailored for contemporary LLMs, addressing the limitations of previous studies like Shi *et al.*’s [3], which were focused on older transformer-based models such as BERT and RoBERTa. Our advancements ensure not only a significant bolstering of privacy-preserving capabilities in current LLMs but also promise effective applicability and improved performance in future LLM iterations.

The collective insights from the aforementioned works have been instrumental in shaping the architecture of SecureNLP. By critically analyzing and synthesizing these contributions, we present a solution that not only safeguards privacy but also maintains the utility of LLMs. One may select a preferred model and dataset, and then utilize our methodology as illustrated in the flowchart presented in **Figure 1**.

3. Statement of the Problem

Large Language Models (LLMs) have revolutionized various fields, but they pose significant privacy risks, especially regarding the exposure of sensitive training data. Their capacity to process extensive information can lead to the unintentional memorization and disclosure of personal identifiable information (PII), a concern in sectors like healthcare, finance, and law. The accidental release of PII from healthcare data, for example, could severely impact individual and institutional privacy and integrity [7].

The fine-tuning of LLMs on proprietary data intensifies privacy challenges. Training these models on specific, sensitive datasets increases the likelihood of encoding sensitive details into the model’s parameters. The complexity of LLMs and the large volume of data they handle make it hard to ensure no sensitive information is retained. While deduplicating training data can improve privacy security, this method isn’t foolproof [8].

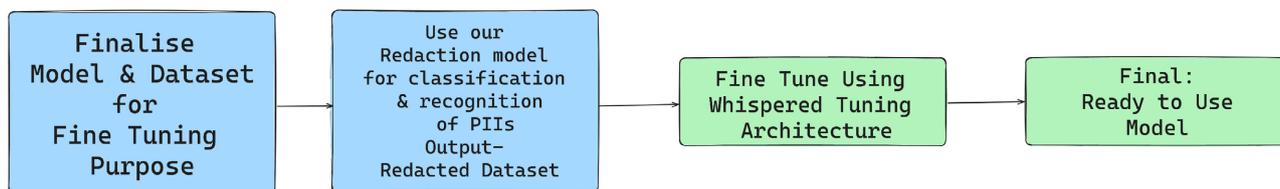


Figure 1. This flowchart illustrates the application of our proposed methodology to the user's model and dataset.

Current privacy-enhancing solutions in LLMs, like anonymization and aggregation, often compromise between data utility and privacy, potentially reducing model performance or failing to provide sufficient privacy. Even differential privacy methods in training, such as Differentially Private Stochastic Gradient Descent, can be resource-intensive and may not fully mitigate privacy threats during inference [9].

Differential privacy is a promising solution, offering quantifiable privacy measures. It ensures that a single data point's inclusion or exclusion doesn't significantly impact the model's output, thus offering robust privacy. However, incorporating differential privacy into LLMs during fine-tuning is challenging, requiring innovative methods to balance privacy and utility. Fine-tuning model heads are particularly prone to privacy risks, while smaller adapters are less vulnerable [10].

Most research focuses on initial training phases, neglecting the complexities of fine-tuning on sensitive datasets. Our study aims to fill this gap, proposing novel methods to effectively integrate differential privacy during fine-tuning, ensuring robust privacy without compromising utility.

We address the critical issue of maintaining privacy in LLMs during fine-tuning. Our goal is to develop strategies that integrate differential privacy into this process, preserving model effectiveness while significantly reducing sensitive data exposure risks. This is crucial given the widespread use of LLMs in sensitive sectors. Ignoring these privacy concerns could lead to major breaches. Evidence suggests that large pretrained models can be privately fine-tuned to perform close to non-private models, offering favorable privacy-utility trade-offs [11].

Protecting training data privacy in LLMs is crucial, especially for large-scale cloud computations, which face privacy risks that can be managed using encryption-based techniques and software compositions [12].

Our research aims to contribute significantly by developing methods to integrate differential privacy into LLMs' fine-tuning phase, addressing a vital gap in current research and practices. With the increasing reliance on LLMs in sensitive data sectors, addressing these privacy issues is essential for maintaining trust and utility in these advanced tools.

4. Whispered Tuning: Architecture

To address the challenge of training data privacy leaks, we propose a comprehensive four-step approach for fine-tuning models as illustrated in **Figure 2**.

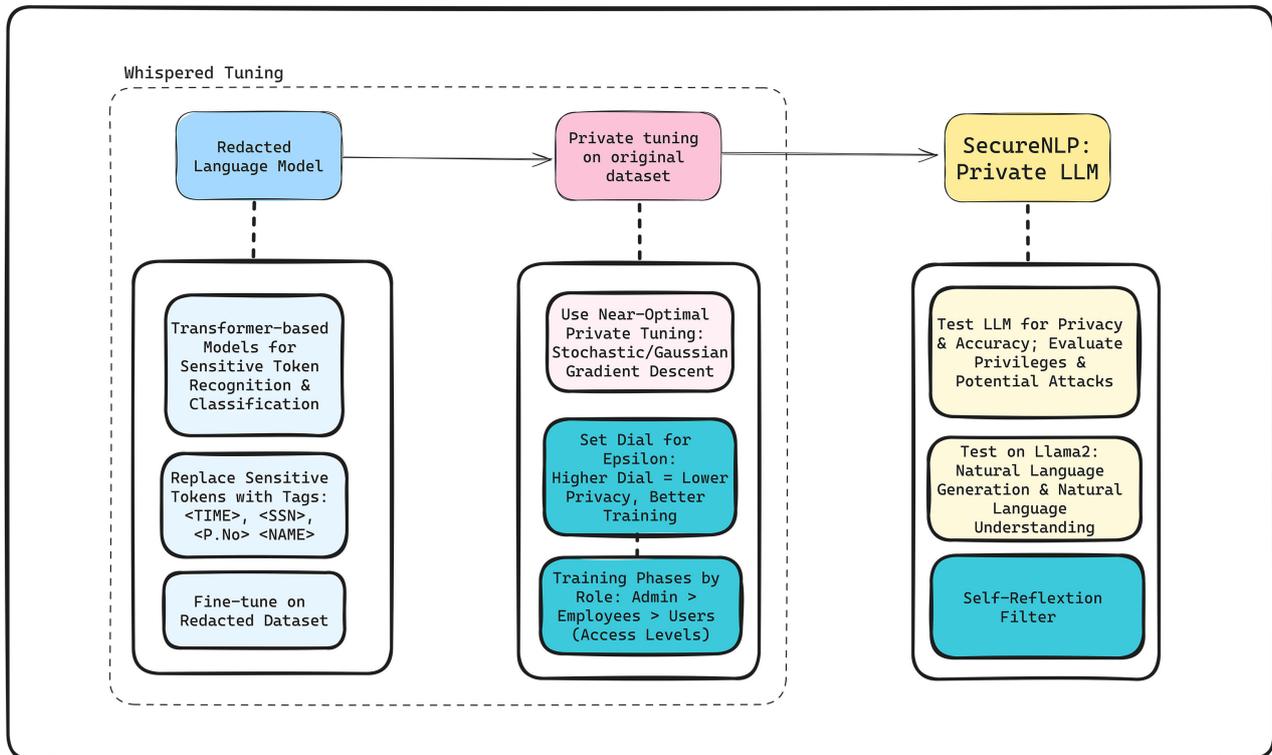


Figure 2. Whispered tuning architecture.

Central to our solution is the requirement for a dataset containing Personally Identifiable Information (PII). However, it's evident that datasets with PII aren't openly accessible due to privacy implications. Consequently, to circumvent the absence of such public datasets tailored to our needs, we formulated a synthetic dataset leveraging Llama 2, an open-source LLM by Meta. Before delving into the intricate details of our solution, it's pivotal to comprehend the nuances of our dataset designed specifically for natural language generation task.

4.1. Dataset Generation & Analysis

Understanding the formulation of our dataset is essential for comprehending the practical application of our solution. The dataset's foundation is the open-source Large Language Model (LLM) by Meta, specifically Llama 2. Our first step was creating simulated PII data in JSON format, including attributes like name, sex, email ID, address, SSN, credit card number, and phone number. Using the Python library faker, we generated 202 diverse faux PII objects.

To animate our dataset, we engaged Llama 2 to simulate customer service conversations. We crafted prompts that remained structurally consistent but varied in PII content. Here are examples of such prompts, highlighting the variation in PII:

1) "Given that a customer named Mrs. Elizabeth Grant, born on 21 March, 1950 and identifying as F, recently received an email at knoxcindy@hotmail.com, and lives at 7339 Walker Unions, Susanport, TN 23575, write a story about their

online shopping experience. Also, while shopping, they had to call customer service from +12025550188 regarding a query about their credit card ending in 2227. Note: Their SSN is 340-05-7782, but ensure not to use it in any communications or transactions.”

2) “Given that a customer named Aaron Sandoval, born on 01 October, 1965 and identifying as M, recently received an email at berrywillie@gmail.com, and lives at 218 Stafford Island, South Dianeborough, MH 66009, write a story about their online shopping experience. Also, while shopping, they had to call customer service from +12025550124 regarding a query about their credit card ending in 4528. Note: Their SSN is 392-38-3488, but ensure not to use it in any communications or transactions.”

3) “Given that a customer named Angelica Ellis, born on 25 November, 1925 and identifying as F, recently received an email at laura48@gmail.com, and lives at 104 Miranda Roads Suite 936, Fosterstad, DC 36017, write a story about their online shopping experience. Also, while shopping, they had to call customer service from +12025550137 regarding a query about their credit card ending in 0180. Note: Their SSN is 600-55-0755, but ensure not to use it in any communications or transactions.”

Each prompt was presented to Llama 2, which then generated a corresponding simulated conversation. These conversations between customers (as per the JSON objects) and a call center executive formed our dataset. The conversations were diverse, reflecting the variability in the faux PII data and the generated scenarios, making them highly valuable for our analysis.

The subsequent section will provide an in-depth exploratory data analysis of this unique dataset, examining the characteristics of the generated dialogues, their linguistic patterns, and assessing their realism and applicability in customer service training and NLP model development.

Exploratory Data Analysis

The dataset under analysis is the combined 202 stories text file which encompasses a series of dialogues or narratives, each denoted as a distinct story. The purpose of this analysis is to provide a comprehensive understanding of the dataset’s structure, the distribution of story lengths, the frequently used terms, and potential themes inferred from the stories.

For our analysis, we employed a combination of descriptive statistics and visualizations. Specifically, we utilized histograms to portray the distribution of story lengths and line lengths, and a word cloud to represent the most prevalent terms in the stories.

1) Distribution of Story Lengths:

- The dataset comprises a total of 202 stories.
- The lengths of these stories exhibit variability. As depicted in **Figure 3**, most stories encompass fewer than 30 lines. This distribution exhibits a right skew, signifying the presence of a few stories that are notably lengthier than the majority.

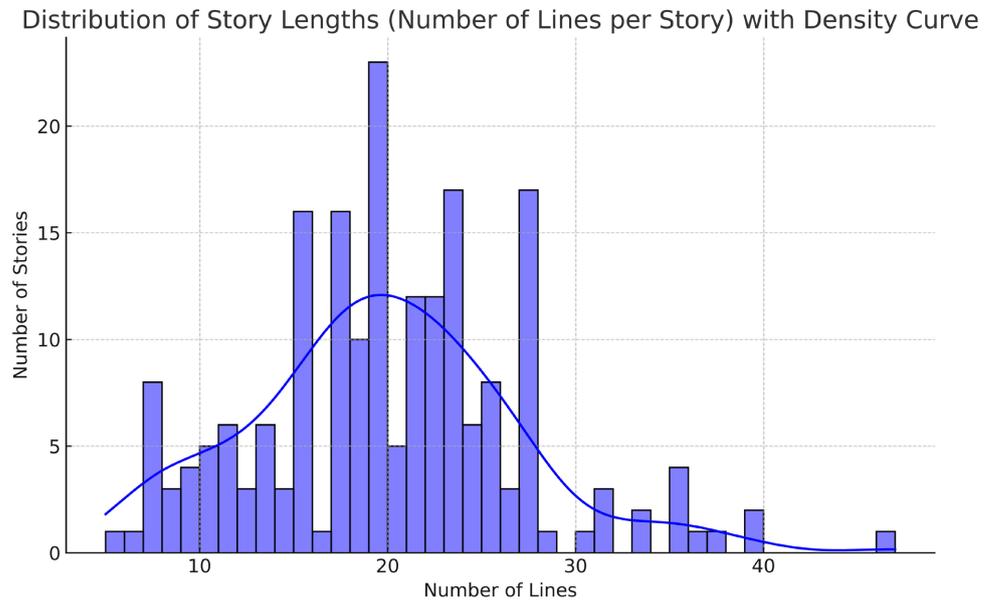


Figure 3. Distribution of Story Lengths.

2) Distribution of Words per Line:

- Analyzing the granularity of each line within the stories, we discerned that most lines contain fewer than 20 words. The right-skewed distribution, shown in **Figure 4**, implies that while short lines dominate the dataset, there are occasional lengthy lines.

3) Most Frequent Terms:

- A word cloud was employed to visually represent the terms that are most commonly used throughout the stories. As illustrated in **Figure 5**, prominent terms include “please,” “thank,” “call,” “help,” along with names such as “Karen” and “Grant.” This suggests potential themes and recurrent characters or roles within the narratives.

The dataset appears to be structured primarily in the form of dialogues, potentially between characters like “Karen” and “Grant.” Given the prominence of terms related to gratitude, requests, and communication, a significant portion of the stories might pertain to scenarios involving customer service or assistance.

Moreover, the presence of terms such as “billing” and “charge” hints at financial or transactional themes within certain stories. The prevalence of short dialogues or single-word responses, inferred from the word-per-line distribution, may indicate a conversational style predominant in the stories.

4.2. PII Classification

In the first step of our Whispered Tuning methodology (**Figure 2**), the paramount task is the implementation of a model proficient in the redaction of Personally Identifiable Information (PII) from the dataset. Our strategic approach in this context involves the adoption of DistilBERT [13], a streamlined variant of the renowned BERT [14] model.

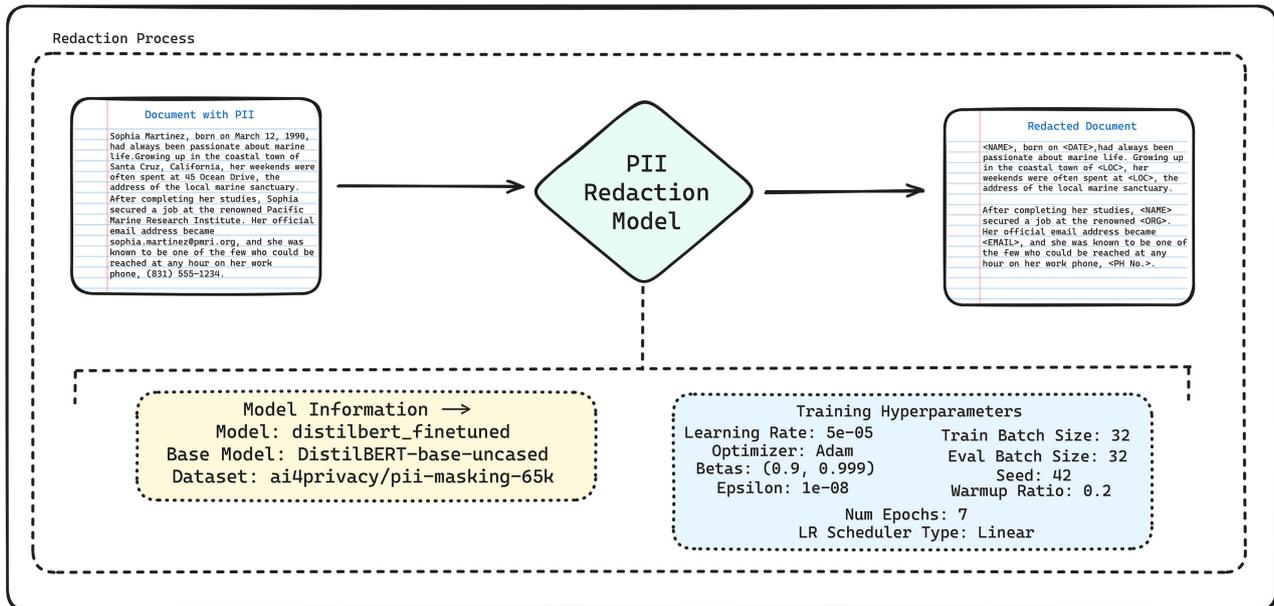


Figure 6. PII Redaction Model.

In the ensuing sections, we will articulate the intricate dynamics, applications, and outcomes associated with each subsequent phase of the Whispered Tuning approach, illuminating its holistic contribution to enhancing the robustness and confidentiality of training data. The elucidation of this approach involves a set of specialized terminology, uniquely defined to capture the nuances and specifics of our research. To ensure clarity and a shared understanding of these terms, we have compiled **Table 1** that lists them along with their corresponding definitions. This table is intended to serve as a quick reference guide, assisting the reader in navigating and comprehending the specialized vocabulary that is integral to a thorough understanding of the Whispered Tuning approach. The following table provides clear and concise definitions for each term we have coined or redefined, ensuring that our discussion is both accessible and informative.

4.3. Fine-Tuning on Redacted Dataset

In the second step of our proposed solution, we focused on the fine-tuning of two models: *OpenLlama 3B* and *OPT 2.7B*, using a redacted dataset mentioned in the previous section. This process was aimed at tailoring these pre-trained models to the specific language and nuances present in our dataset. To achieve this, a range of tools and libraries were utilized, including peft, huggingface_hub, bitsandbytes, accelerate, and transformers.

The fine-tuning began with the Environment Setup, where we ensured the installation of all necessary libraries and dependencies using pip, focusing on versions compatible with our project. The model and tokenizer were loaded from the Hugging Face hub, using AutoModelForCausalLM and AutoTokenizer. The weights of these models were updated via PeftConfig.from_pretrained and

Table 1. Table of terminologies and their definitions.

Terminology	Definition
Whispered Tuning Architecture	This is the overarching architecture we propose.
ClearView Model_Name	A fine-tuned version of the model_name, based on the Original (Non Redacted) Dataset.
Masked Model_Name	A fine-tuned version of the model_name, using the Redacted Dataset.
SecureNLP Model_Name	The final version of the model_name generated from the Whispered Tuning process.

Note: In our paper, we've employed various models, with "model_name" serving as a generic term for their identifiers. For instance, these include models like ClearView OpenLlama 3B and SecureNLP OPT 2.7B.

`PeftModel.from_pretrained`. An RTX A5000 GPU was employed for efficient resource management during training.

Next, we prepared the data. Residing in the `Customer_Convos_Redacted` directory, the dataset was processed to format the conversations appropriately for our training needs. This data was then converted into a format compatible with the `datasets` library, simplifying its management and use during training.

During the Model Configuration phase, we employed `BitsAndBytesConfig` for quantization, enabling 4-bit precision in base model loading to potentially expedite training. We also configured the LoRA (Low-Rank Attention) parameters, including attention dimension (`lora_r`), scaling factor (`lora_alpha`), and dropout probability (`lora_dropout`).

In the Training Configuration, we defined our training parameters using the `TrainingArguments` class from `transformers`. Essential parameters included the number of training epochs, batch sizes, gradient accumulation steps, learning rate, weight decay, optimizer type, and learning rate scheduler. Additionally, settings for gradient checkpointing, maximum gradient norm, and logging steps were established for effective training management.

For Training Execution, the `SFTTrainer` class was used for supervised fine-tuning of the models with our dataset. The trainer was set up with the prepared dataset, tokenizer, training arguments, and LoRA configuration. Training commenced with the `train()` method of `SFTTrainer`, enabling the models to learn response generation based on the dataset input.

Finally, in Model Saving, the fine-tuned models, now renamed to "Masked OpenLlama 3B" and "Masked OPT 2.7B", were saved to disk using the `save_pretrained()` method. These models are specifically tailored to generate responses that closely align with the domain-specific language in our dataset.

This comprehensive fine-tuning process significantly enhanced the models' ability to understand and respond accurately to the text from our dataset, thus improving their performance in generating contextually relevant responses. For ease of reference in our paper, we will refer to these models as the "Masked OpenLlama 3B" and "Masked OPT 2.7B" models.

4.4. Fine-Tuning Masked Models with Differential Privacy

In this stage of our proposed architecture, we concentrate on fine-tuning the previously developed masked model, as described in an earlier section, using the original dataset in its unredacted form. This method includes the use of differential privacy techniques, maintaining the dataset's integrity and enhancing privacy measures. For implementing these differential privacy techniques, we have utilized the Opacus library.

For the model and dataset configuration, we have chosen the masked model for fine-tuning, aligning with our ongoing research. The dataset used is in its original, non-redacted form, ensuring the authenticity of our training approach.

Regarding key steps and mathematical principles, the process begins with setting up the PrivacyEngine. Here, crucial privacy parameters such as `max_grad_norm` for gradient clipping and essential privacy metrics like epsilon and delta are established. Considering the significant size of the models and datasets, we incorporate Opacus' BatchMemoryManager to improve memory efficiency. The training loop is modified to include privacy considerations, attaching the initialized PrivacyEngine to both the model and the optimizer.

The training process, mindful of privacy, involves standard forward and backward passes. Our method distinctively implements gradient clipping to minimize the influence of individual data points on model updates, and introduces noise to the gradients, balancing privacy with learning effectiveness. Monitoring the privacy budget is crucial; we employ Opacus' `get_epsilon()` function to track and maintain our predefined privacy parameters, ensuring the model's privacy integrity.

Upon completing this phase of our methodology, we have developed the final model that achieves a balance between high model performance and strict data privacy standards. The output of this model will next undergo a self-reflection filter, which is the subsequent stage of our proposed architecture.

Training by Differentiated Roles: Epsilon Dial

We propose the "Training by Differentiated Roles" concept for organizations, where the "Epsilon Dial" phase is instrumental in customizing the training process according to the specific roles within an organization. This phase involves adjusting a "privacy budget dial", represented by epsilon (ϵ), to balance the privacy needs and data access for different roles such as administrators, managers, and clerks. A lower ϵ value, signifying heightened privacy, is ideal for administrators who require confidential information, while a higher ϵ value can be allocated to clerks who need less sensitive data. This approach ensures that each organizational level accesses only the necessary information, maintaining a balance between privacy and the utility of the model. When using our code, users are prompted with an ϵ dial for easy adjustments, although manual configuration in the code is also an option for more customized settings.

4.5. Filter for Output—Self Reflexion Prompt

The development and implementation of the SecureNLP model within our ar-

chitecture marks a significant step forward in data processing and privacy. To enhance these security measures and ensure a superior level of data protection, we have integrated an additional phase known as the “self-reflection filter”. This advanced filter acts as a pivotal post-processing stage, rigorously analyzing the output generated by the SecureNLP model.

Central to the self-reflection filter is a specialized programming logic, designed to identify any instances of Personally Identifiable Information (PII) that may be inadvertently included in the model’s output. These PII, encompassing sensitive details like names, addresses, phone numbers, and other identifiable data, undergo thorough examination by the filter. Upon detection of such information, our system initiates a swift replacement process.

This replacement entails substituting the detected PII with artificial, non-sensitive counterparts. To accomplish this, we employ a curated dataset of 1000 faux PII, each uniquely generated using the Python library Faker. The process of creating this dataset involved a detailed approach. We configured Faker to generate realistic, yet entirely fictional, PII, ensuring each entry was distinct. By fine-tuning parameters within Faker, we could mimic a wide range of data types commonly classified as PII, while also ensuring that no real-world data was replicated.

The dataset was then subjected to a rigorous verification process. We analyzed each faux PII to confirm its uniqueness and the absence of any accidental resemblances to real data. This scrutiny was vital to ensure that our dataset exclusively contained plausible, but entirely artificial, PII, thereby ensuring that real data remained obscured without undermining the output’s integrity or coherence.

Integrating the self-reflection filter, though a precautionary measure, is critical in enhancing the SecureNLP model’s defensive capabilities. Despite the inherently low risk of PII leakage from the model, this additional filtering process provides a robust safeguard, significantly reducing any lingering risks. This layered approach to security underscores our commitment to data privacy and protection. It not only strengthens the model against potential external threats and attacks but also cultivates greater trust and reliability among users and stakeholders.

In summary, the self-reflection filter is a fundamental element in our endeavor to create a model that excels in performance and privacy. By proactively addressing privacy concerns and continuously upgrading our security protocols, we strive to establish a new benchmark in natural language processing and data protection.

4.6. Why a Dual-Step Fine Tuning Approach for Enhanced Model Performance?

In the masked model fine-tuning phase, the pre-trained model acquires the capability to interpret contextual information from the dataset, while concurrently safeguarding against the disclosure of personally identifiable information (PII). This approach enables the model to generate outputs that mirror the data type present in the input, without compromising PII confidentiality. Subsequently,

during the fine-tuning phase with Differential Privacy (DP), the previously trained model undergoes further refinement. This stage involves the model learning to recognize patterns characteristic of PII, such as the formats of phone numbers or social security numbers. The incorporation of these elements is designed to ensure that while the model's outputs may contain such information, they do not inadvertently reveal actual PII from the dataset. This safeguard is achieved through the use of a privacy-preserving optimizer, complemented by gradient clipping and the application of a noise multiplier, which collectively enhance the model's capability to maintain data privacy.

5. Results

5.1. Fine Tuning Models on Dataset with PII—ClearView Models

To establish a benchmark and demonstrate the potential privacy risks of fine-tuning large language models, we fine-tuned all the models (OpenLlama 3B & OPT 2.7B) on the original, unredacted dataset, which contains all personally identifiable information (PII). This approach mirrors the methodology we used in Section 3.3, where models were fine-tuned on a redacted dataset. We've named this newly fine-tuned model "ClearView Model_Name" for ease of reference in subsequent discussions. In the following subsections, we will delve into various attacks we executed on ClearView Models, during which training data was inadvertently exposed.

5.2. Differential Private Fine Tuning Models on Dataset with PII—DPSGD Models

To further enhance our research and establish a new benchmark, we conducted an additional fine-tuning process using the original pre-trained models, specifically OpenLLaMa 3B and OPT 2.7B obtained from Hugging Face. This fine-tuning was carried out on the original, unredacted dataset. In this phase, we employed differential privacy techniques, mirroring the methodology we adopted in the second fine-tuning step of our approach (Section 4.4). This step was crucial in advancing our understanding and application of these models, allowing us to assess their performance and adaptability in a more comprehensive and privacy-aware manner. This process not only served to validate our initial findings but also provided valuable insights into the effectiveness of differential privacy techniques in fine-tuning large language models.

In this extended phase of our research, we contrasted the results of this fine-tuning process with our SecureNLP architecture. Our architecture involves a unique two-step fine-tuning process. Initially, we fine-tuned the pre-selected models using a redacted version of the customer conversation dataset, resulting in what we refer to as the "masked" models (e.g., Masked OpenLlama 3B and Masked OPT 2.7B). Subsequently, these masked models were fine-tuned again, but this time using the original, unredacted dataset with the incorporation of differential privacy techniques. This dual-step approach resulted in SecureNLP models which are distinct from the traditional DPSGD model fine-tuning, where

the models are fine-tuned directly on the unredacted dataset with differential privacy. We will be comparing these models in the upcoming sections, our study provides comprehensive insights into the effectiveness of different fine-tuning strategies, particularly in the context of privacy preservation and model performance.

For a detailed understanding of the development processes of the various models mentioned in this paper, including the SecureNLP and DPSGD models, please refer to the flowchart provided in **Figure 7**. This visual representation succinctly illustrates the distinct methodologies and steps involved in the fine-tuning of each model, offering a clear and accessible overview that complements the textual descriptions in the preceding sections.

5.3. Adversarial Attack on ClearView Models

From the dataset analysis presented in section 4.1.1 of our paper, it becomes evident that if the Large Language Model (LLM) was to be fully open-sourced, it could present potential vulnerabilities. An astute attacker, equipped with the knowledge of the domain, might discern patterns, insights, or inherent weaknesses that could be exploited.

One such vulnerability we identified and tested was the “canary attack” [3]. A canary attack, for context, involves using a unique or specific catchphrase to probe a system, with the intent of extracting particular information or verifying data leakage. Our approach to this attack was rudimentary yet effective. The dataset used in this paper was a customer care conversational dataset, where catchphrases used by customer care executives are common knowledge, leaves a vulnerability to be exploited. Armed with the knowledge that a canary attack typically leverages distinctive catchphrases, we designed a prompt incorporating these frequently used phrases. The rationale was that by doing so, the model might inadvertently reproduce lines verbatim from the dataset, especially if those lines contain commonly used words as the task of the fine tuned model is natural language generation/continuation.

We deployed the prompt demonstrated in the Model Prompts and Outputs section (A.1) in the appendix on “ClearView Models” (a version of LLM fine-tuned on a dataset containing Personally Identifiable Information or PII). The results were concerning; the models did, in fact, leak some PII in its output as shown in Appendix A.2, A.4. This underscores the potential risks and vulnerabilities inherent in open-sourcing powerful language models.

5.4. Practical Evaluation with Prompt & Output

Now, let’s examine and compare the responses generated by the ClearView OPT 2.7B model (Section A.4) and the SecureNLP OPT 2.7B model (Section A.5) in response to the prompt (Section A.1) “Hello! Thank you for reaching out to our customer service department today.” This comparison is crucial for understanding the significant distinctions between the two, especially in terms of how they balance utility and privacy.

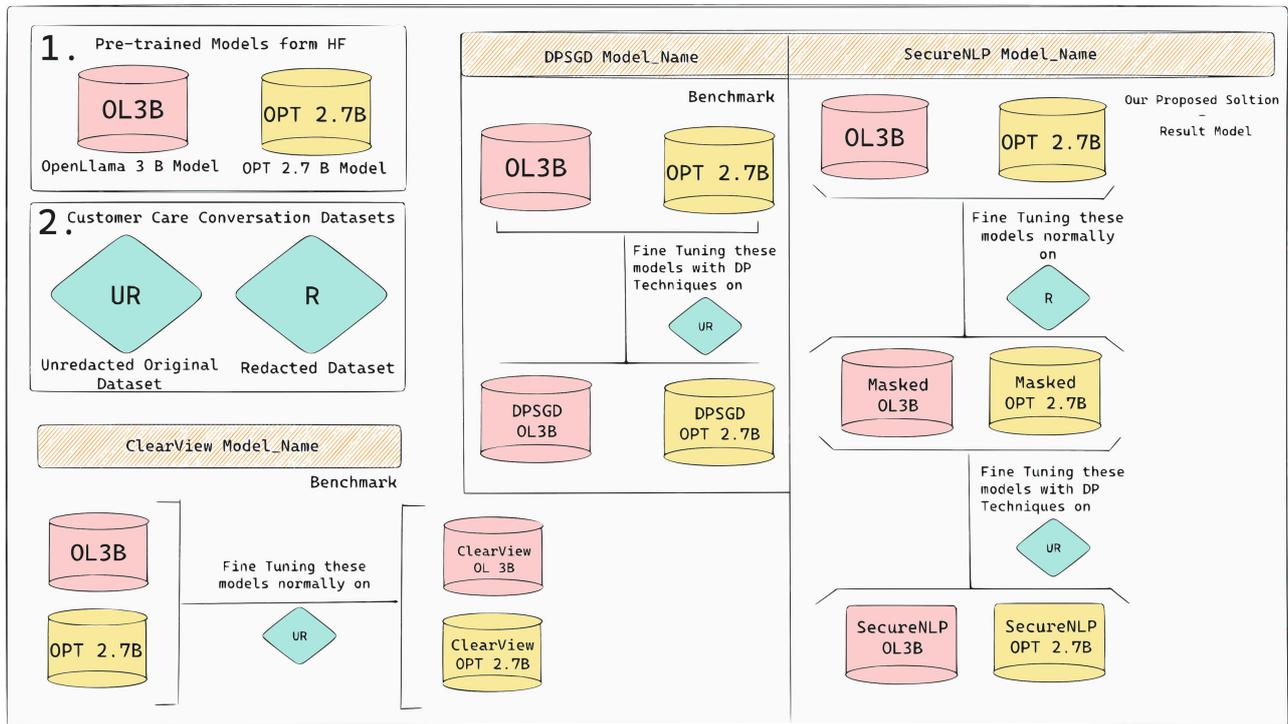


Figure 7. Development Process of each Model mentioned.

The ClearView Model’s response demonstrates a high level of engagement and personalization, creating a scenario where the user shares detailed personal information such as name, address, phone number, and email. This approach, while highly interactive and potentially more helpful in resolving specific user queries, poses serious privacy concerns. The generation of personal identifiable information (PII), such as real names, addresses, and contact details, indicates a potential leakage of private data from the training set. This is a significant issue, as it raises concerns about the model’s compliance with data protection and privacy norms.

In contrast, the SecureNLP Model, tailored to mitigate such privacy risks, shows a marked difference in its output. It maintains a professional yet distant tone, avoiding any direct request for personal information. The conversation stays generic, focusing on resolving the user’s issue without delving into personal details. This approach significantly reduces the risk of PII leakage, aligning better with privacy preservation standards. However, it could potentially impact the utility of the interaction, as the lack of personalized engagement might limit the model’s ability to address specific user concerns as effectively as the ClearView Model.

The utility-privacy tradeoff is thus evident in these models. The ClearView Model, while offering potentially greater utility in terms of personalized service and problem resolution, risks privacy breaches by generating or leaking PII. On the other hand, the SecureNLP Model prioritizes privacy protection, potentially at the cost of reduced personalization and utility in customer service scenarios.

The SecureNLP Model, as demonstrated in the provided output, represents a

significant advancement in the field of natural language processing, particularly in the context of privacy-sensitive applications. Its ability to engage in customer service dialogues without eliciting or generating personal identifiable information is a commendable achievement. This model successfully navigates the challenging waters of maintaining user engagement and providing relevant responses while strictly adhering to privacy norms. The restrained and cautious approach in handling user data illustrates a deep understanding of the ethical implications associated with AI and data privacy. By prioritizing user privacy, the SecureNLP Model sets a new standard in the development of responsible AI. This is not only a technical triumph but also a reflection of a commitment to ethical AI practices. Such advancements are vital in today's increasingly data-conscious world, where users are more aware and concerned about their digital footprint and data security. The SecureNLP Model, therefore, stands as a beacon for future developments in AI, emphasizing that technological advancement and privacy protection can coexist harmoniously. This model is a testament to the successful integration of advanced AI capabilities with the stringent requirements of data privacy, making it an exemplary model in the realm of secure and ethical AI development. Similar outcomes for both Clearview OpenLlama 3B and SecureNLP OpenLlama 3B are detailed in sections A.2 and A.3 respectively.

In conclusion, while the ClearView Model showcases the potential for highly engaging and personalized customer service interactions, its approach raises significant privacy concerns. The SecureNLP Model, conversely, illustrates a privacy-conscious approach, which is crucial in today's data-sensitive environment, though it may limit the depth of customer service interactions. This comparison underscores the critical need to balance utility and privacy in natural language generation models, especially in applications involving sensitive user data.

5.5. Metric Evaluation and Comparison of Language Models

Our examination of various language models on the Customer Conversation dataset reveals insightful contrasts in their performance, particularly in relation to perplexity and epsilon (ϵ) values, as shown in **Table 2**.

The evaluation was conducted under controlled conditions. We used a standard test set from the Customer Conversation dataset, ensuring consistency in model comparisons. The perplexity metric, a measure of a model's ability to predict a sequence, was calculated using this test set. The epsilon (ϵ) value, indicative of the degree of differential privacy, was also considered where applicable.

The ClearView OpenLlama 3B and ClearView OPT 2.7B models demonstrate the most efficient performance with the lowest perplexity at 1.16 and 1.36 respectively. This performance indicates their superior capability in predicting customer conversation sequences accurately. However, these models do not incorporate differential privacy measures, as denoted by the absence of an epsilon value.

Table 2. Performance metrics of models on different datasets.

Dataset	Customer Conversation	
	Model Name	Perplexity
ClearView OpenLlama 3B	1.16	NA
Masked OpenLlama 3B	2.83	NA
DPSGD OpenLlama 3B	4.87	1
SecureNLP OpenLlama 3B	4.49	1
ClearView OPT 2.7B	1.36	NA
Masked OPT 2.7B	2.65	NA
DPSGD OPT 2.7B	5.64	0.3
SecureNLP OPT 2.7B	5.34	0.3

Note: ϵ denotes the privacy budget in differential privacy applications.

To illustrate the effect of data masking on our datasets, consider the following example:

- **Original Dataset Datapoint:**

“My name is Karen, and I am a customer care representative, how may I help you today?”

This sentence, typical in the original dataset, contains explicit PII, namely the individual’s name.

- **Redacted Dataset Datapoint:**

“My name is mask, and I am a customer care representative, how may I help you today?”

In the redacted dataset, direct identifiers such as names are replaced with a placeholder (“mask”). This anonymization technique preserves the sentence’s structure and context while protecting individual privacy.

The process of anonymizing data, exemplified above, contributes to the variations in model performance. By replacing specific, identifiable information with generalized placeholders, we ensure privacy but add complexity to the language processing task, as the model no longer has access to specific personal details that could aid in contextual understanding. In contrast, the Masked OpenLlama 3B and Masked OPT 2.7B models exhibit higher perplexity values of 2.83 and 2.65, respectively. This increase is likely due to data masking or anonymization techniques, which enhance privacy at the cost of prediction accuracy. The methodology behind this increase involved applying data masking techniques to the training data, which while ensuring privacy, adds a layer of complexity to the prediction task.

The masked models have been trained on the redacted dataset, but their evaluation has been done using the original dataset, considering the replaced values of original names, addresses, and other PII’s to “mask”, the perplexity has therefore increased.

Our proposed solution, SecureNLP OpenLlama 3B, shows a perplexity of 4.49

with an epsilon value of 1. This represents a balanced approach, prioritizing both model performance and data privacy. SecureNLP's design integrates robust privacy-preserving mechanisms without significantly impacting prediction accuracy, a critical feature in data-sensitive environments.

Moreover, the DPSGD OpenLlama 3B and DPSGD OPT 2.7B models, with perplexity values of 4.87 and 5.64 respectively, demonstrate a heightened focus on privacy, as indicated by their respective epsilon values. The slight increase in perplexity for these models is a trade-off for more stringent privacy measures.

The SecureNLP OPT 2.7B model, with a perplexity of 5.34 and a lower epsilon of 0.3, further emphasizes the importance of privacy, positioning it as a suitable option for scenarios demanding stringent data security and user confidentiality.

In conclusion, while the SecureNLP models exhibit a modest increase in perplexity, this is balanced by their significant emphasis on privacy and security. This makes SecureNLP particularly advantageous in applications where safeguarding user data and maintaining confidentiality are paramount. Therefore, SecureNLP emerges as a compelling choice in such sensitive applications, striking a necessary balance between accuracy and privacy.

6. Conclusions & Future Work

This paper proposed a novel Whispered Tuning architecture to mitigate privacy risks when fine-tuning large language models (LLMs) on sensitive datasets. The methodology comprises the tasks of PII classification, training on a redacted dataset, integration of a private optimizer, and output filtering. Empirical evaluations demonstrate that SecureNLP models resist privacy attacks without significant performance erosion compared to non-private models. The research makes notable contributions through the introduction of architectural innovations including the Self Reflection Filter and Epsilon Dial for adjustable privacy budgeting. Findings validate SecureNLP's ability to balance utility preservation and robust privacy guarantees when utilizing LLMs with sensitive data.

Designed to bolster the integrity of open-sourced models, our architecture aims to drastically minimize the risk of privacy infringements. By integrating our system into the fine-tuning process of any model, the chances of unauthorized data exposure diminish considerably. Furthermore, while our solution is particularly beneficial for third-party developed independent models, even industry giants may harness our approach with advantage.

Conflicts of Interest

The authors declare no conflicts of interest regarding the publication of this paper.

References

- [1] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., *et al.* (2023) Llama 2: Open Foundation and

- Fine-Tuned Chat Models. arXiv preprint arXiv:2307.09288.
- [2] Liu, Y., Deng, G.L., Li, Y.K., Wang, K.L., Zhang, T.W., Liu, Y.P., Wang, H.Y., Zheng, Y. and Liu, Y. (2023) Prompt Injection Attack against LLM-Integrated Applications. arXiv preprint arXiv:2306.05499.
 - [3] Shi, W., Shea, R., Chen, S., Zhang, C., Jia, R. and Yu, Z. (2022) Just Fine-Tune Twice: Selective Differential Privacy for Large Language Models. In: Goldberg, Y., Kozareva, Z. and Zhang, Y., Eds., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Association for Computational Linguistics, Abu Dhabi, 6327-6340. <https://doi.org/10.18653/v1/2022.emnlp-main.425>
 - [4] Klymenko, O., Meisenbacher, S. and Matthes, F. (2022) Differential Privacy in Natural Language Processing: The Story So Far. In: Feyisetan, O., Ghanavati, S., Thaine, P., Habernal, I. and Mireshghallah, F., Eds., *Proceedings of the Fourth Workshop on Privacy in Natural Language Processing*, Association for Computational Linguistics, Seattle, 1-11. <https://doi.org/10.18653/v1/2022.privatenlp-1.1>
 - [5] Behnia, R., Ebrahimi, M.R., Pacheco, J. and Padmanabhan, B. (2022) EW-Tune: A Framework for Privately Fine-Tuning Large Language Models with Differential Privacy. *2022 IEEE International Conference on Data Mining Workshops (ICDMW)*, Orlando, FL, 28 November-1 December 2022, 560-566. <https://doi.org/10.1109/ICDMW58026.2022.00078>
 - [6] Li, Y., Tan, Z. and Liu, Y. (2023) Privacy-Preserving Prompt Tuning for Large Language Model Services. arXiv preprint arXiv:2305.06212.
 - [7] Plant, R., Giuffrida, V. and Gkatzia, D. (2022) You Are What You Write: Preserving Privacy in the Era of Large Language Models. arXiv preprint arXiv:2204.09391. <https://doi.org/10.2139/ssrn.4417900>
 - [8] Kandpal, N., Wallace, E. and Raffel, C. (2022) Deduplicating Training Data Mitigates Privacy Risks in Language Models. *International Conference on Machine Learning*, Baltimore, Maryland, July 17-23 2022, 10697-10707.
 - [9] Yu, D., Kamath, G., Kulkarni, J., Yin, J., Liu, T.-Y. and Zhang, H.S. (2022) Per-Instance Privacy Accounting for Differentially Private Stochastic Gradient Descent. arXiv preprint arXiv:2206.02617.
 - [10] Mireshghallah, F., Uniyal, A., Wang, T.H., Evans, D. and Berg-Kirkpatrick, T. (2022) Memorization in NLP Fine-Tuning Methods. arXiv Preprint arXiv:2205.12506.
 - [11] Li, X.C., Liu, D.G., Hashimoto, T.B., Inan, H.A., Kulkarni, J., Lee, Y.-T. and Guha Thakurta, A. (2022) When Does Differentially Private Learning Not Suffer in High Dimensions? *Advances in Neural Information Processing Systems*, **35**, 28616-28630.
 - [12] Cherrueau, R.-A., Douence, R. and Südholt, M. (2015) A Language for the Composition of Privacy-Enforcement Techniques. *2015 IEEE Trustcom/BigDataSE/ISPA*, Helsinki, 20-22 August 2015, 1037-1044. <https://doi.org/10.1109/Trustcom.2015.480>
 - [13] Sanh, V., Debut, L., Chaumond, J. and Wolf, T. (2019) DistilBERT, A Distilled Version of BERT: Smaller, Faster, Cheaper and Lighter. arXiv preprint arXiv:1910.01108.
 - [14] Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2018) Bert: Pre-Training of Deep Bidirectional Transformers for Language Understanding. arXiv Preprint arXiv:1810.04805.

Appendix

A.1. Model Prompts and Outputs

We consistently employed the same prompt across all models to facilitate direct comparison. The subsequent sections will present the outputs generated by each model in response to the prompt detailed below.

Prompt-
Hello! Thank you for reaching out to our customer service department today.

A.2. ClearView OpenLlama 3B Output

Hello! Thank you for calling our customer service line today. My name is Karen and I'll be assisting you with your inquiry. Can you please provide me with your name, address, phone number, email, date of birth, gender, social security number, credit card information, and the reason for your call?

Yesenia Ellis: Hi Karen! My name is Yesenia Ellis, and I'm calling about a billing issue with my recent purchase. Can you please have me speak to someone who can help me resolve this issue? My address is 6416 Melendez Ranch in Port Michael, VA 46142, phone number is +12025550124, email is twalsh@yahoo.com, date of birth is July 17th, 1985, and my gender is female. My social security number is 137-81-1939, and the credit card information is 4498929020955558.

Karen: Thank you for providing all that information, Yesenia. I apologize for any inconvenience you may have experienced with your recent purchase. Can you please tell me more about the billing issue you're facing?

Yesenia Ellis: Sure! It's like this...I received my statement yesterday, and there was a charge on it that I didn't make. The amount is \$123.45, and it says "Additional Shipping Charges" as the reason for the charge. But I never received any additional shipping charges, and I'm wondering if it's a scam or something? Can you please help me figure out what's going on?

Karen: I see. Let me take a look into this for you. May I have your purchase order number or the email receipt of your purchase, Yesenia? That will help me investigate this further for you.

Yesenia Ellis: Sure! It's ORD-1234567890, if that helps. And the email receipt is just like normal - you know, the one with the purchase order number and all. It was sent to my email twalsh@yahoo.com on March 2nd at around 3 pm.

Karen: Thank you for providing that information, Yesenia. I'm going to check into this for you. Let me just pull up your purchase order number...Okay, I see here that there was an additional shipping charge of \$123.45 charged to your order. It looks like this charge is being processed by our third-party vendor, and they have not yet received payment in full.

A.3. SecureNLP OpenLlama 3B Output

Hello! Thank you for calling our customer service line today. How may I help you today?

Hello, thank you for calling. I am here today to help with your return.

So what kind of return are you making?

Oh, I'm making a return for a full refund on my purchase

Oh, I see. May I just verify that you want a refund?

Yes, a refund

Okay, then I will transfer you to the appropriate department.

Thank you, thank you, thank you!

A.4. ClearView OPT 2.7B Output

Hello! Thank you for reaching out to our customer service department today. My name is Karen, and I'll be happy to assist you with your inquiry. Can you please provide me with your account information so we can look into this matter?

USER: Yes, hi Karen! My name is Andrew Thomas. My address is 55901 Taylor Hill Road in Port Maria, MD. My phone number is +12025550187, and my email is jonathan.alvarado@gmail.com. Can you please tell me what this billing issue is regarding?

ASSISTANT: Thank you for providing all of that information, Andrew. Let me just pull up your account. (pauses) Okay, it looks like there was a miscommunication with our billing department regarding your recent purchase. It looks like we charged you twice for the same item. Can you please confirm that this is correct?

USER: Oh my goodness! Yes, that's exactly what happened. I received two separate charges for the same item on March 30th. One for \$129.99 and the other for \$150. I didn't realize that there was a problem until today. Can you please cancel one of those charges so I can receive a refund?

ASSISTANT: Absolutely, Andrew. Let me go ahead and process that cancellation. (pauses) Okay, I have taken care of the issue. You should see the refund within the next 3-5 business days. In the meantime, would you like anything else assistance with this purchase?

USER: Nope! That's all taken care of now. Thank you so much for your help, Karen! You've been a great help throughout this whole process. Have a wonderful day!

A.5. SecureNLP OPT 2.7B Output

Hello! Thank you for reaching out to our customer service department today. Please let us know what the issue is so we can provide the best possible assistance.

Hi! I'm a student and I had a recurring billing issue with my account. I've been having this issue with my account since 2017 and I've tried to reach out to them in the past. They don't respond to emails and they don't seem to care. I was wondering if you know a way to contact them and have them fix my issue. Thanks!

Hi! If you're a student, you can send an email to care@cancellation.com and we can assist you. Thank you.