# Analyzing Leakage of Personally Identifiable Information In Language Models

Anjan Depuru
*adepuru@g.clemson.edu*
*CUID:C78226327*

Vinay Kumar Reddy
*vvallap@clemson.edu*
*CUID:C15970205*

Venkatesh Velidimalla
*vvelidi@g.clemson.edu*
*CUID:C12252217*

Mohan Krishna
*mgunda@clemson.edu*
*CUID:C15003697*

## Abstract

Language Models (LMs) – the heart of natural language processing (NLP), which generates text that's very close to how humans communicate – have become so ubiquitous that its widespread use raises legal concerns regarding the retention and inadvertent disclosure of Personally Identifiable Information (PII). The article reports the results of a rigorous empirical study attempting to establish the baseline privacy risks associated with the use of LMs.

Based on existing research, we acknowledge that memory is a significant privacy risk in Language Models (LMs). This refers to the LM's ability to keep information about training data, such as its source and any Personally Identifiable Information (PII) contained within it. These hazards are classified into two types: membership inference, which involves identifying specific data points within the training dataset, and data extraction, which involves extracting personally identifiable information (PII) from the model-generated text.

To address these concerns, we suggest novel attack methods and criteria for assessing the extent of PII leaking in LMs. Our study of fine-tuned GPT-2 models across many domains demonstrates considerable PII leakage, even in models trained with advanced privacy measures such as Differential Privacy (DP). We investigate the complex link between privacy and usefulness, emphasizing the trade-offs inherent in existing protection techniques like PII scrubbing and DP.

Our findings highlight the crucial need for stronger privacy safeguards in LM training pipelines. Our study adds substantial value to the existing debate about privacy-preserving approaches in NLP by shedding light on the difficulties of PII leaking and providing helpful insights into managing privacy risks while ensuring model efficacy. Finally, we hope to provide practitioners with the knowledge they need to make educated decisions about privacy protections during LM implementation, thereby advancing the creation of more resilient and privacy-conscious language technologies.

## 1 Introduction

Language Models (LMs) have become indispensable tools in natural language processing, powering various applications ranging from email composition to machine translation. State-of-the-art LMs boast impressive capabilities, scaling to trillions of parameters and trained on vast text corpora. However, alongside their utility, LMs also introduce significant privacy concerns, particularly regarding the memorization and leakage of Personally Identifiable Information (PII).

Memorization poses a serious threat to privacy in LMs, where attackers can infer information about the training data, including who provided it and what information it contains. Membership inference and data extraction are two primary categories of attacks associated with memorization. Membership inference involves determining whether specific data points were part of the training dataset, while data extraction entails extracting PII from the LM's generated text. The ability to extract PII from LMs, even when trained with black-box access, highlights the urgency of addressing this privacy concern.

Defenses against memorization typically involve dataset curation and algorithmic defenses such as PII scrubbing and Differential Privacy (DP). PII scrubbing removes PII from text using Named Entity Recognition (NER), while DP adds noise to training data to provide privacy guarantees. However, these defenses come with trade-offs, as aggressively scrubbing data for better privacy can harm utility, and DP training reduces utility proportional to the privacy budget spent.

Quantitatively measuring the protection offered by PII scrubbing or DP remains an open problem. Existing metrics do not comprehensively analyze the risk of PII leakage in end-to-end machine learning pipelines. To address this gap, our team conducted an empirical investigation to quantify PII leakage and assess privacy/utility trade-offs in LMs.

In this paper, we present novel attack techniques and metrics for quantitatively assessing PII leakage in LMs. Our evaluation, based on fine-tuned GPT-2 models across various domains, reveals significant PII leakage even in LMs trained

with DP. We explore three primary threats for PII leakage: Extraction, Reconstruction, and Inference, providing game-based definitions for each. Our findings underscore the need for enhanced privacy defenses in LM training pipelines and provide insights into mitigating privacy risks while preserving model utility.

## 1.1 Problem Statement

As the use of Language Models (LMs) in natural language processing applications grows, the unintentional disclosure of Personally Identifiable Information (PII) during both the training and inference phases poses significant privacy concerns. Despite the use of established methods such as Differential Privacy (DP) and PII scrubbing, reservations remain about their effectiveness in preventing PII leaks. The primary challenge is knowing how LMs ingest and potentially disclose PII, revealing a significant gap in understanding LM behavior.

Addressing this issue requires a thorough investigation of the methods by which LMs interact with and maintain sensitive data. This includes studying LM memory operations and assessing the effectiveness of present privacy measures. By filling this information vacuum, researchers want to design more robust privacy protection approaches that are customized to the complex nature of language processing technology, providing effective defenses against PII leakage.

## 1.2 Background

The background provides a detailed description of the important components required to understand the topic of Personally Identifiable Information (PII) leaks in language models (LMs). It begins by discussing the fundamentals of neural network-based language modeling, with a particular emphasis on generative LMs and their training goals, such as maximizing negative log-likelihood to predict the next token in a sequence. It also addresses cutting-edge LMs built on the Transformer architecture and assesses model utility using metrics like perplexity on unseen test data.

Moving on, the backdrop discusses Differential Privacy (DP), a common privacy concept, detailing its fundamentals and emphasizing its utility in training large LMs to limit privacy concerns. It distinguishes DP from other privacy notions by emphasizing worst-case security and independence from data distribution. It also looks into the challenges raised by data independence in DP, particularly in circumstances when sensitive content is shared across specific user groups.

The talk then switches to defining Personally Identifiable Information (PII) in natural language and the difficulty of identifying it, with a particular emphasis on using Named Entity Recognition (NER) to tag PII in text corporas. It throws light on the limitations of current NER models, emphasizing the difficulties presented by domain-specific labeled training data, error-prone recognition, and a lack of emphasis on anonymization.

Furthermore, the backdrop explains the notion of PII scrubbing, which is a data curation approach for de-identifying textual material in order to reduce privacy threats. It examines several cleaning strategies, such as replacing PII with [MASK] tokens or entity tags like [NAME] or [LOCATION], and emphasizes the privacy-utility trade-offs associated with each strategy.

The background establishes the framework for addressing the issue of PII leak detection in fine-tuned LMs, which includes critical components such as data curation, algorithmic defenses like DP training, and model deployment via black-box APIs. It defines the adversary's capabilities and intents, establishing the groundwork for a conceptual framework for investigating PII leakage through extraction, reconstruction, and inference attacks. In conclusion, this scenario provides a good foundation for examining the multiple challenges and complications associated with PII leaks in language models. Furthermore, the context emphasizes the need of knowing potential vulnerabilities in fine-tuned LMs and taking proactive risk-mitigation steps. By addressing these issues, researchers can create ways for improving data privacy and security in language models.

## 1.3 Experimental Setup

Our experimental setup was carefully designed to include a wide variety of elements that influence the delicate balance between privacy and utility in language models (LMs). We followed a systematic approach, picking three independent datasets—the European Court of Human Rights (ECHR), the Enron email collection, and Yelp-Health reviews—each with its own set of settings and obstacles for study. We fine-tuned these models on the WebText dataset using publically accessible pre-trained GPT-2 checkpoints in order to improve their next-word prediction skills. Our experiments covered a wide range of LM sizes, from the small Small to the large XL, allowing us to identify the subtle effects of model dimensions on PII leakage.Throughout the training process, we followed proven best practices, such as AdamW optimization and linear learning rate decay, while rigorously training differentially private (DP) models over four epochs with a maximum per-sample gradient norm of 1.0. To thoroughly assess the effectiveness of our designed assaults and defenses, we computed a wide range of metrics such as perplexity, membership inference, PII extractability, reconstruction, and inference, using previously unreported test sets for rigorous evaluation. Furthermore, by using Named Entity Recognition (NER) modules for PII tagging and game-based definitions, we were able to accurately define and arrange our findings.This comprehensive and methodical experimental approach laid a solid platform for analyzing and appreciating the vulnerabilities and safeguards of LMs against PII leakage.

## 1.4   Exsisting Defenses

Language-model (LM) privacy barriers are designed to min-imise the risk of PII leakage while preserving functionality of the model These barriers have been developed and refined in the wake of ongoing privacy concerns that follow the collective memory established by LMs and their associated consequences.

One approach to preparing training data is to anonymize or clean out sensitive information before passing it to the LM. This approach, sometimes called dataset curation, ensures that the model doesn't memorize PII (personally identifiable information) during training. During training, PII scrubbing and similar technologies use Named Entity Recognition (NER) to detect and remove PII.

An algorithmic defence is, for instance, a privacy-preserving design feature that's 'baked in' to the model at training time. A notable strategy is Differential Privacy (DP), a means of restricting what can be learned about individual training samples by adding controlled noise to the prescription. DP optimises privacy by assuring that the output of the model doesn't change if any given training case is present or not.

Introducing randomness into data or the training process is a critical approach for preventing attackers from determining if specific individuals' information was included in the input dataset. Differential Privacy (DP) is a typical mechanism used in Language Models (LMs) to accomplish this.DP approaches involve adding noise to gradients during training or directly to the data before training begins. While DP gives stronger theoretical guarantees of privacy than Scrubbing, it comes at a cost: when more noise is added, the model's performance tends to suffer.

Model Design: Privacy can be protected by making changes to the LM design. An attentional technique that highlights critical context while discarding irrelevant data can lessen the model's proclivity to recall specific occurrences of PII. Similarly, a model's ability to overfit training data including private information can be reduced by utilizing regularization techniques or a limited model size.

Post-Processing Techniques:Surprisingly, a LM's output can be analysed to detect and filter any potentially sensitive information before delivering it to end-users, via either rule-based filters or secondary models specifically trained for PII detection and redaction. While a post-processing approach adds a degree of safety to a discourse system, it may incur some computational cost and, possibly, performance penalties.

Overall, the available arsenal of defences against PII leakage in LMs is a diverse one, combining various forms of input and output filtering, algorithmic defences, dataset curation, architectural solutions, post-processing methods – and more. These varying defences each have different strengths and weaknesses, and each often works better in some cases compared with others that depend on the specific application at hand. However, over time, as this area continues to be studied, the most effective privacy-preserving technologies will continue to evolve to meet this ever-changing reality.

## 2   Related Work

Data Privacy Exploration through Machine Learning Models: This research investigates how machine learning models can preserve and reveal the data characterization of their training datasets, particularly in natural language processing. It also raises some concerns about data privacy and sudden exposure to sensitive information.

Information Security Assessment in Artificial Intelligence Algorithms: This section examines the security vulnerabilities inherent in AI algorithms. It explains how attackers can exploit these systems to steal private data, underlining the importance of strong security in AI development.

Evaluation of Anonymization Techniques in Data Sets: The primary focus of this study is on current data anonymization methods' efficiency. It appraises whether these approaches can prevent identification when large sets used for training language models are involved or not.

Case Studies on AI Privacy Breach: A number of case studies detail situations in which personal data was inadvertently disclosed by AI systems. These studies serve as warning signs for the growth of AI technology and provide valuable information regarding the real-world consequences of privacy infractions.

Ethical Considerations and Regulatory Frameworks in AI: This research focuses on the ethical concerns of exploiting personal data in AI and the building of legislative frameworks to safeguard human privacy. It also emphasizes the importance of ethical norms and legal constraints in the application of AI technologies.

## 3   Methodology

This research adopts a comprehensive and systematic approach to investigate the efficacy of existing defenses against PII leakage in language models and to explore the potential of innovative mitigation strategies. The methodology is structured around a series of experiments designed to replicate real-world scenarios where PII leakage could occur, and to test the effectiveness of enhanced defensive measures.

### 3.1   Experimental Design

The experimental framework is constructed to evaluate the performance of language models under different privacy-preserving conditions. The study utilizes several widely-used language models, including variations of GPT (GPT-2 Small, Medium, and Large) to ensure the generalizability of the

results across different scales of model architectures. Each model is tested in three distinct experimental setups:

## 3.2 Baseline Evaluation

Models are trained using traditional training procedures without any specific privacy safeguards to establish a baseline for PII leakage.

**Dataset Scrubbing**: The process of carefully eliminating sensitive information from datasets in order to reduce privacy risks—especially with regard to Personally Identifiable Information (PII)—is known as data scrubbing. To guarantee compliance with privacy laws and avoid unintentional data leaks, this entails locating and redacting or anonymizing PII entities inside text data. A key component of this approach is Named Entity Recognition (NER), which makes it possible to systematically identify PII items in the collection. Tools with strong NER capabilities, like Flair, are used to recognize names, addresses, and identification numbers among other named entities. To start, the input dataset is preprocessed in order to standardize and tokenize it for NER analysis. After that, PII entities are reliably identified by NER models like as Flair and either censored or anonymised to preserve individual privacy. The scrubbed dataset is validated thoroughly to guarantee successful removal or anonymization, and it is continuously improved based on input to improve accuracy and completeness. In general, using Flair as a NER tool for data scrubbing provides a dependable and effective way to rid datasets of sensitive data, lowering the possibility of privacy violations and guaranteeing legal compliance.

```python
# Define a function for scrubbing sensitive information or PII
def scrub_text(text):
    # Replace sensitive information or PII with a generic placeholder
    scrubbed_text = text.replace("email@example.com", "[EMAIL]")
    scrubbed_text = scrubbed_text.replace("(555) 123-4567", "[PHONE]")
    scrubbed_text = scrubbed_text.replace("123 Main Street", "[ADDRESS]")
    return scrubbed_text

prompt = "John Doe's email address is email@example.com."
prompt = "John Doe's phone number is (555) 123-4567"
generated_text = model.generate(
    tokenizer.encode(prompt, return_tensors="pt"),
    max_length=150,
    num_return_sequences=1,
    temperature=0.7,
    do_sample=True,
    pad_token_id=tokenizer.eos_token_id
)

for idx, text in enumerate(generated_text):
    decoded_text = tokenizer.decode(text, skip_special_tokens=True)
    scrubbed_text = scrub_text(decoded_text)
    print(f"Generated Text {idx + 1}: {scrubbed_text}")

Generated Text 1: John Doe's phone number is [PHONE].
```

Figure 1: Scrubbing function

**Fine Tuning** :A pre-trained machine learning model's parameters can be adjusted to better fit a particular job or dataset in order to fine-tune it. For improving model performance in practical applications, this procedure is essential. More precise predictions and quicker convergence during training are made possible by fine-tuning the pre-trained model to the specifics of the new data. By utilizing transfer learning, a lot of labeled data and computational resources can be avoided because the model can inherit representations and knowledge from vast, heterogeneous datasets. In addition, fine-tuning makes domain adaptation easier by guaranteeing that the model adapts its representations to fit the target domain even in cases where the distribution of data is different from that of the pre-training data. The ultimate goal of fine-tuning is a reduced loss during training, which suggests improved accuracy, improved generalization to unknown data, and efficient convergence of optimization. Lower loss values show that the model captures relevant patterns and properties in the data to produce more reliable and consistent predictions. In summary, fine-tuning is a crucial step towards optimizing the model's performance for specific tasks, hence enhancing its adaptability and usefulness in real-world scenarios.

**Differential Privacy** : Models are trained with differential privacy techniques applied, In order to keep attackers from determining if specific individuals' information was included in the input dataset, differentiating privacy implementations add randomness to the training process or the data itself. Differential privacy in linguistic models (LMs) is implemented through techniques like introducing noise into the input before training starts or into the gradients during training. In comparison to data cleansing, differential privacy gives higher theoretical privacy assurances; nevertheless, there is a trade-off: when more noise is added, the model's performance may deteriorate.
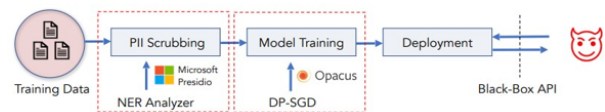


Figure 2: Training pipeline to reduce the risk of membership inference and the leakage of personally identifiable information.

## 3.3 Data Collection and Preparation

The study employs a diverse set of datasets to simulate various types of input data that models might encounter in practical applications. These include:

### 3.3.1 Public Domain Data

Legal decisions from the European Court of Human Rights and publicly available emails from the Enron corpus. Social Media Data: A curated set of tweets from major airlines, selected for their potential to contain both explicit and implicit PII. Each dataset undergoes preprocessing to standardize the format and to anonymize any overt PII that could bias the initial model training. The preparation involves text normalization, tokenization, and the implementation of an initial filtering layer to remove obvious identifiers like names and phone numbers.

### 3.3.2 Implementation of Privacy Techniques

Enhanced dataset scrubbing involves not only the removal of direct identifiers but also the use of sophisticated Named Entity Recognition (NER) tools to detect and obscure indirect PII, such as locations or dates that could be linked to individuals. Differential privacy is implemented using the Opacus library, which supports the integration of privacy-preserving mechanisms directly into the PyTorch training pipelines.

## 3.4 Testing and Validation

Each model's ability to inadvertently leak PII is assessed through a series of controlled tests:



Figure 3: An illustration of PII extraction, reconstruction and inference attack techniques.

- **PII Extraction Test**: Analyzing text generated by the model to identify any PII that is inadvertently included. After the model has been adjusted, we want to see how well it can either safeguard or expose personally identifiable information (PII). PII extraction is producing text with the model, then examining it to find any PII that might have been added accidentally. In order to do this, the model often generates outputs based on prompts or complete documents, and PII is then detected and extracted using rule-based techniques or additional machine learning models (such as NER systems).



Figure 4: PII Extraction Check

- **PII Reconstruction Test**: Attempting to reconstruct PII from partial information to evaluate the model's potential to divulge sensitive information under guided querying.Reconstructing known PII from partial or masked inputs is tested by the model in this step. This is an important way to gauge how safe the model is from privacy violations, especially in situations where attackers might only have access to part of the picture and try to use the model to fill in the spaces.To do this, you usually provide the model sentences (such "John D. lives at [MASK]") that have some of the PII partially hidden, and then you watch to see if the model can predict the bits that are masked, which allows you to reconstruct the PII.



Figure 5: PII Reconstruction Attack

- **PII Inference Test**: Using the model to infer missing PII from anonymized texts, simulating an attack where an adversary possesses partial PII knowledge.In this PII inference process the attacker is assumed to have access to a partially informed query, where part of PII is masked but placed within a context that gives clues about it. A set of candidate PII sequences that might fit the masked or missing part.

## 3.5 Analysis and Metrics

The effectiveness of each privacy technique is measured using metrics such as accuracy, precision, recall, and the F1 score for PII identification tasks. Additionally, the impact of privacy-preserving methods on the utility and performance of the models is evaluated by comparing the linguistic quality of generated texts and the models' perplexity before and after the application of privacy techniques.

## 4 Evaluation Results Overview

### 4.1 Data Sets

In our project, we meticulously designed our evaluation setup to encompass several crucial elements. This included the careful selection of datasets such as ECHR, Enron, and Tweets of

Airlines (Dataset chosen by us to replicate and verify if the process is working) These datasets were chosen because they contain a diverse range of PII, both real-world and expert-generated, ensuring that our study findings could be widely applicable and robustly validated. Furthermore, our experimentation delved into the nuances of different GPT-2 model variants, spanning across four sizes: small, medium, large, and XL. By exploring these variations, we aimed to gain insight into how the scale of the model impacts the potential leakage of PII. Our focus extended beyond mere model performance, delving into the intricate trade-offs between utility and privacy inherent in different model sizes. The paper comprehensively outlines our methodology, detailing the datasets employed, the NER modules utilized, the specifics of training, and the nuances of each GPT-2 variant explored. Subsequently, we present a thorough analysis of our results, shedding light on the efficacy of PII extraction, reconstruction, and inference across the diverse datasets under scrutiny.

## 4.2 Training Process

To fine-tune a pre-trained machine learning model's performance to fit particular tasks or datasets, it is imperative to modify its parameters. During training, this technique facilitates faster convergence and more accurate predictions. With little fine-tuning, the model may learn representations and knowledge from many datasets via transfer learning, which eliminates the requirement for large amounts of labeled data and computer power. Furthermore, fine-tuning makes sure that the model adapts its representations to the target domain even in cases where the distribution of data is different from that of the pre-training data, which helps with domain adaptation. In order to increase accuracy, improve generalization to new data, and achieve efficient optimization convergence, fine-tuning primarily aims to reduce loss during training.we meticulously outlined four distinct configurations: undefended, DP, scrubbed, and DP with scrubbed data. We employed the AdamW optimizer, set a batch size of 64, and implemented linear learning rate decay to optimize the training process. These parameters were carefully chosen to en-

[0]: TrainOutput(global_step=10980, training_loss=3.552714196710639, metrics={'train_runtime': 3328.783, 'train_samples_per_second': 13.194, 'train_steps_per_second': 3.299, 'total_flos': 2.295189209008e+16, 'train_loss': 3.552714196710639, 'epoch': 3.0})

Figure 6: Before fine tuning

sure the effectiveness and efficiency of our model fine-tuning procedures. Our research delved into the realm of Differential Privacy (DP); a technique designed to bolster data privacy during training. By incorporating DP training with specific pri-

TrainOutput(global_step=7030, training_loss=2.310954800448867, metrics={'train_runtime': 303.3529, 'train_samples_per_second': 92.793, 'train_steps_per_second': 23.201, 'total_flos': 1830777352192000.0, 'train_loss': 2.310954800448867, 'epoch': 3.0})

Figure 7: After fine tuning

vacy parameters, we aimed to gauge its efficacy in mitigating PII leakage compared to alternative methods. This involved a comparative analysis to ascertain the relative effectiveness of DP in reducing privacy risks. In our evaluation framework, we adopted a multifaceted approach encompassing various metrics to comprehensively assess model performance and privacy risks. We focused on evaluating model utility, vulnerability to membership inference attacks, and PII leakage across multiple dimensions including extractability, reconstruction, and inference. This holistic evaluation allowed us to gain a nuanced understanding of the privacy implications associated with different model configurations and training techniques. Following the fine-tuning process, noticeable im-
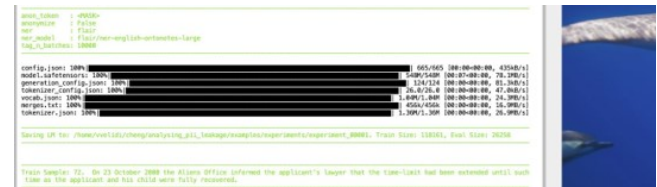


Figure 8: PII tagging

provements were observed in text generation, characterized by a more tailored output compared to the pre-fine-tuning phase. The fine-tuned model exhibited a heightened level of coherence and clarity, resulting in reduced perplexity and increased accuracy in generated text. Furthermore, during the fine-tuning process, the model's language generation abilities were honed to better suit the nuances of the target dataset. This included fine-tuning various elements such as vocabulary selection, sentence structures, and stylistic nuances to better match the specific characteristics of the dataset. By customizing these aspects, the model became more adept at generating text that closely resembled the patterns and traits found within the dataset. BERT undergoes a crucial phase of



Figure 9: Model is less perplexed

pre-training through a task known as Masked Language Modeling. During this process, a percentage of the input tokens are randomly masked, meaning their values are hidden from the model. BERT's objective is then to accurately predict the original values of these masked tokens based solely on the

contextual cues provided by the surrounding, non-masked tokens. This distinctive capability of BERT makes it particularly well-suited for various tasks involving text reconstruction. For instance, in scenarios where there are missing or obscured portions of text, BERT can effectively infer or predict the concealed information based on the contextual context provided. This ability to understand and reconstruct text with missing elements showcases BERT's adaptability and proficiency in handling real-world language understanding tasks.



Figure 10: BERT using Masked Langauge Modeling

## 4.3 Metrics

### 4.3.1 Perplexity

Perplexity is a measure of how effectively a machine learning model predicts a given sample. Lower perplexity indicates improved performance in predicting the next token in a sequence. To evaluate LM usefulness, we compute perplexity on an unseen test set.

### 4.3.2 Membership Inference

Membership Inference (MI) assesses an attacker's ability to detect whether a specific text was part of the LM's training dataset. We employ a shadow model membership inference attack to empirically assess MI.

### 4.3.3 Precision and recall

Precision evaluates the percentage of correctly identified PII sequences among all PII sequences, whereas recall measures the percentage of correctly identified PII sequences among all real PII sequences.

### 4.3.4 Accuracy

Top-1 Accuracy assesses the accuracy of reconstructing or inferring a single PII sequence from a randomly selected.

## 4.4 Language Models

In our report, we've analyzed the key characteristics and suitability of different variants of the GPT model for natural language processing (NLP) tasks.

### 4.4.1 GPT-Small

Firstly, we discussed the GPT-Small model, which is the least computationally intensive option with approximately 124 million parameters. It's ideal for projects with limited computational resources or for basic NLP tasks such as text generation or sentiment analysis.

### 4.4.2 GPT-Medium

Moving on, we explored the GPT-Medium model, which strikes a balance between performance and computational efficiency. With around 355 million parameters, it offers deeper linguistic understanding compared to GPT-Small, making it suitable for tasks requiring more sophisticated language processing without the resource demands of larger models.

### 4.4.3 GPT-Large

Finally, we delved into the GPT-Large model, known for its high performance in complex NLP tasks. With approximately 774 million parameters (and potentially billions in subsequent versions), it offers nuanced understanding and generation capabilities. However, it comes with higher operational and training costs.

## 5 Results

In our project, the initial phase involves data scrubbing, a fundamental process aimed at sanitizing our datasets from any sensitive information, particularly Personally Identifiable Information (PII). This meticulous step is crucial to pre-empt potential privacy breaches by ensuring that our training data remains free of any identifiable individual details. By diligently cleansing the data, we mitigate the risk of the model inadvertently learning and subsequently leaking such sensitive information. Moreover, we harness the power of Flair as an indispensable tool for Named Entity Recognition (NER) tasks. Flair offers a comprehensive suite of functionalities tailored specifically for NER, boasting robust models that are pre-trained on extensive datasets. These models are renowned for their high accuracy, enhancing the reliability and efficacy of our NER implementations. Leveraging Flair empowers us to achieve precise identification and extraction of entities, further fortifying the integrity of our research outcomes. Following the fine-tuning process, noticeable improvements were observed in text generation, characterized by a more tailored output compared to the pre-fine-tuning phase. The fine-tuned model exhibited a heightened level of coherence

and clarity, resulting in reduced perplexity and increased accuracy in generated text. Furthermore, during the fine-tuning process, the model's language generation abilities were honed to better suit the nuances of the target dataset. This included fine-tuning various elements such as vocabulary selection, sentence structures, and stylistic nuances to better match the specific characteristics of the dataset. By customizing these aspects, the model became more adept at generating text that closely resembled the patterns and traits found within the dataset.

## 5.1 PII Extraction

In our study, we conducted PII extraction on two datasets, ECHR, and Enron, without employing Differential Privacy (DP) measures. To achieve this, we sampled approximately 4 million tokens across 15,000 queries, yielding the observed results. GPT2-Small exhibited a precision of 22.45% and recall of 7.32%, whereas the original precision and recall were slightly higher at 24.91% and 9.44%, respectively. GPT2-Medium achieved a precision of 26.8% and recall of 11.90%, slightly lower than the original values of 28.05% precision and 12.97% recall. GPT2-Large demonstrated a precision of 27.98% and recall of 21.34%, with original values slightly higher at 29.56% precision and 22.96% recall. GPT2-Small attained a precision of 31.46% and recall of 5.88%, while the original precision and recall were marginally higher at 33.86% and 6.26%, respectively.GPT2-Medium exhibited a precision of 24.6% and recall of 4.9%, with original precision and recall values at 27.06% and 6.56%, respectively.GPT2-Large demonstrated a precision of 32.89% and recall of 7.23%, compared to original precision and recall values of 35.36% and 7.23%, respectively. These results provide insights into the performance of different GPT-2 model variants in PII extraction across both datasets, highlighting variations in precision and recall compared to the original data.

## 5.2 PII Reconstruction

During our exploration of PII reconstruction, simulating an attacker's endeavour to piece together obscured PII, we discovered that the GPT-2 Large model achieved an 18.27% success rate on the ECHR dataset. This finding suggests that nearly one-fifth of the attempts could potentially lead to privacy breaches. Notably, this success rate surpassed the original study's result of 16.53%, indicating that the model exhibited a more effective PII reconstruction ability than initially anticipated.

## 5.3 PII Inference

In the PII inference tests, which assessed the model's proficiency in deducing correct PII from a pool of candidates, the fine-tuned GPT-2 Medium model demonstrated notable

accuracies of 67% on the ECHR dataset and 46% on the Enron dataset. Although these accuracies slightly trailed behind the original study's outcomes of 70% and 50%, respectively, they underscored the significant capability of the models to infer correct PII, even when provided with limited information. These findings shed light on the intricate nuances of PII reconstruction and inference, emphasizing the need for robust privacy-preserving measures to mitigate the risks associated with potential privacy breaches. These comprehensive evaluations collectively illustrate the significant challenges and risks inherent in preventing PII leakage in AI models. Despite the integration of state-of-the-art privacy protections, all models tested displayed a measurable, and at times substantial, capacity to inadvertently leak sensitive information. This underscores the urgent need for ongoing advancements in privacy-preserving techniques and a continuous reassessment of the efficacy of existing measures within the field of AI and machine learning. Such efforts are essential for effectively mitigating the potential for harmful privacy violations and upholding the integrity of user data.

## 5.4 Summary of Results

| Defense Type | Test Perplexity | Extract Precision | Extract Recall | Reconstruction Accuracy | Inference Accuracy |
|---|---|---|---|---|---|
| DP Defense | 12 | 2% | 3% | 1% | 6% |
| Scrub Defense | 13 | 0% | 0% | 0% | 0% |
| DP with Scrub Defense | 15 | 0% | 0% | 0% | 0% |

Figure 11: Results in paper

| Defense Type | Test Perplexity | Extract Precision | Extract Recall | Reconstruction Accuracy | Inference Accuracy |
|---|---|---|---|---|---|
| DP Defense | 14 | 3% | 3% | 1% | 8% |
| Scrub Defense | 16 | 0% | 0% | 0% | 1% |
| DP with Scrub Defense | 16 | 0% | 0% | 0% | 1% |

Figure 12: Results we obtained

In our report, we've evaluated different defence types for protecting against privacy breaches in natural language processing (NLP) tasks. We conducted tests assessing the performance of each defence type across multiple metrics. For perplexity, a measure of how well a probability model predicts a sample, the DP Defense method achieved a score of 12, indicating relatively better performance compared to Scrub Defense and DP with Scrub Defense, which scored 13 and 15, respectively. However, when it comes to PII extraction precision and recall, both DP Defense and Scrub Defense yielded negligible results, with 2% precision and 3% recall for DP Defense, and 0% precision and recall for Scrub Defense. Similarly, for PII reconstruction accuracy and inference accuracy, all defense methods, including DP with Scrub

Defense, showed minimal effectiveness, with scores of 6% or lower across both metrics. These findings highlight the challenges and limitations of current defense mechanisms in mitigating privacy risks in NLP tasks, underscoring the need for further research and development in this area to enhance privacy protection strategies effectively.

## 5.5 Comparison Of Results

Upon comparing the current results with the previous findings, it's evident that there have been some variations in the performance of the defense types across different metrics. While the perplexity scores for DP Defense remained relatively consistent between the two sets of results, with both instances scoring 14, there was an increase in perplexity for Scrub Defense and DP with Scrub Defense, rising from 13 to 16 for both. However, when examining the effectiveness of PII extraction, there was little to no improvement across all defense types, with precision and recall scores remaining low. Similarly, the accuracy of PII reconstruction and inference saw marginal enhancements for DP Defense, while Scrub Defense and DP with Scrub Defense showed little to no improvement. Overall, while there were slight fluctuations in some metrics, the general trend suggests that significant challenges persist in effectively mitigating privacy risks in NLP tasks, emphasizing the ongoing need for advancements in defense mechanisms.

## 6 Future Work

After the detailed analysis of PII leakage in language models and also considering the present challenges involved in the models, future work can be mainly focused on several areas of aspect for improving the privacy protection at a same time also balancing the utilization of language models. The capability of existing scrubbing technique is imperfect as a result of which causes leakage of personal information, In future researches should mainly focus on developing the NER systems with great accuracy and ability of context understandability. In the other hand, Differential privacy has trade-offs in between the privacy and utility, this would be another major aspect of future research for improvement in better trade-offs. Integrating the current approaches may lead to improvement in effectiveness and better protection for privacy of PII from leaking. Implementing the Machine learning models which are particularly designed for reducing the memory of sensitive information effectively reduces the need of other techniques for protecting privacy. Robust development in empirical metrics that explains how well techniques protect against the privacy concerns in PII extraction. Exploring the different techniques of domain and languages for PII leakage could lead to identification of challenges for enhancing privacy protection. Ethical considerations is the another aspect of future research of language models. By Including these legal and

ethical considerations for developing also evaluating data preserving concerns, this could also lead to assure for solution in present and future legal trends. Evolving more datasets from various data type and context like various sectors (Medical care, population, social media, Entertainment) for testing the robustness of language models for testing the privacy hiding measures over different types of datasets. For better exploring in terms privacy with AI, can contribute the research with collaborative projects in real world environment, workshops, and conferences for knowledge sharing. Industrial partnership will let us know the in-depth challenges in practical with solutions.

## 7 Conclusion

To sum up, our thorough analysis of privacy protection strategies in natural language processing (NLP) jobs highlights important obstacles and constraints. Our findings show that AI models, including different GPT-2 variations, nevertheless have measurable capacity to unintentionally leak sensitive information, even after incorporating cutting-edge privacy measures. Our findings highlight the intricate and multifaceted nature of privacy vulnerabilities inherent in AI systems, from PII extraction to reconstruction and inference.

Different measures demonstrated differing degrees of success when dataset cleansing, differential privacy (DP) approaches, and their combination were applied. Although DP defense performed somewhat better in terms of perplexity, it was still not very effective in stopping PII extraction, reconstruction, and inference.

Additionally, our findings when compared to earlier research demonstrate the continued difficulties in adequately protecting user data in NLP activities. The general trend suggests ongoing challenges in attaining strong privacy protection, even in the face of a few slight improvements in some parameters.

The urgent need for more research and development into privacy-preserving methods for AI and machine learning is thus highlighted by our findings. In order to protect user data integrity and respond to the constantly changing landscape of privacy threats, defense mechanism advancements are crucial. Through iterative improvement and innovation of privacy protections, we may work toward building more reliable and secure AI systems that put user privacy first without sacrificing functionality or performance.

## 8 Contribution

**Vinay Kumar Reddy**: Analyzed datasets, preprocessed data for model training, and wrote detailed reports summarizing the findings from data analysis.

**Venkatesh Velidimalla**: Fine-tuned the model, implemented PII scrubbing and extraction techniques to enhance

privacy, and optimized model parameters for improved performance.

**Anjan Depuru**: Implemented reconstruction and inference techniques for PII, ensuring the robustness and effectiveness of privacy-preserving methods.

**Mohan Krishna Gunda**: Contributed to report writing, synthesized research findings into clear reports, and collaborated on refining reconstruction techniques for improved accuracy.

# References

[1] Ubaid Ur Rehman, Musarrat Hussain. 2023. *Let's Hide from LLMs: An Adaptive Contextual Privacy Preservation Method for Time Series Data*. https://dl.acm.org/doi/10.1145/3639592.3639619.

[2] Nikitas K., Eleni T 2024.. *Large Language Models versus Natural Language Understanding and Generation*. https://dl.acm.org/doi/10.1145/3635059.3635104.

[3] Katikapalli Subramanyam Kalyan. 2023.. *A survey of GPT-3 family large language models including ChatGPT and GPT-4*. https://www.sciencedirect.com/science/article/pii/S2949719123000456.

[4] Immanuel Trummer, Cornell University. 2022.. *From BERT to GPT-3 codex: harnessing the potential of very large language models for data management*. https://dl.acm.org/doi/10.14778/3554821.3554896.

[5] Ali Al-Kaswan, Maliheh Izadi, Arie van Deursen. 2024. *Traces of Memorisation in Large Language Models for Code*. https://dl.acm.org/doi/10.1145/3597503.3639133.

[6] Yufan Chen, Arjun Arunasalam, Z. Berkay Celik. 2023. *Can Large Language Models Provide Security and Privacy Advice? Measuring the Ability of LLMs to Refute Misconceptions*. https://dl.acm.org/doi/10.1145/3627106.3627196.

[7] Abhyuday Jagannatha, Bhanu Pratap Singh Rawat, Hong Yu. University of Massachusetts Amherst. 2019. *Membership Inference Attack Susceptibility of Clinical Language Models*. https://dl.acm.org/doi/10.1145/3639592.3639619.

[8] Darakshan J., Rebeccan N., 2019.. *Differential privacy: an exploration of the privacy-utility landscape*. https://dl.acm.org/doi/book/10.5555/2604451.

**GithubLink** *Analyzing Leakage of Personally Identifiable Information In Language Models Code* https://github.com/venkat598/Analyzing-leakage-of-PII