# Analyzing Leakage of Personally Identifiable Information in Language Models



https://nilslukas.github.io

Nils Lukas, Dec 11 2023
Research Presentation @Meta

**CrySP**
Cryptography, Security, and Privacy
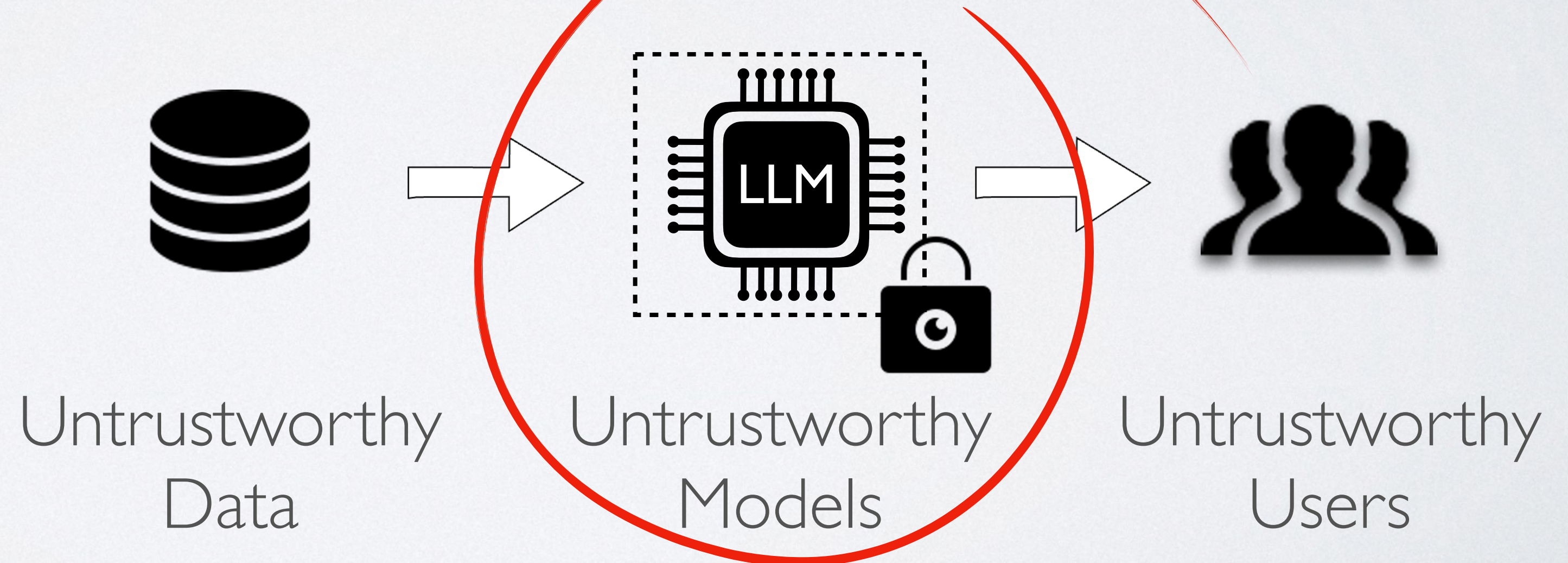Research Group

UNIVERSITY OF
**WATERLOO**

# My Areas of Research



Nils Lukas

**Private Computation**

- Private Set Intersection
- Secure Inference

**Machine Learning**

- Reliability
- Privacy
- Safety

Untrustworthy Data → Untrustworthy Models → Untrustworthy Users
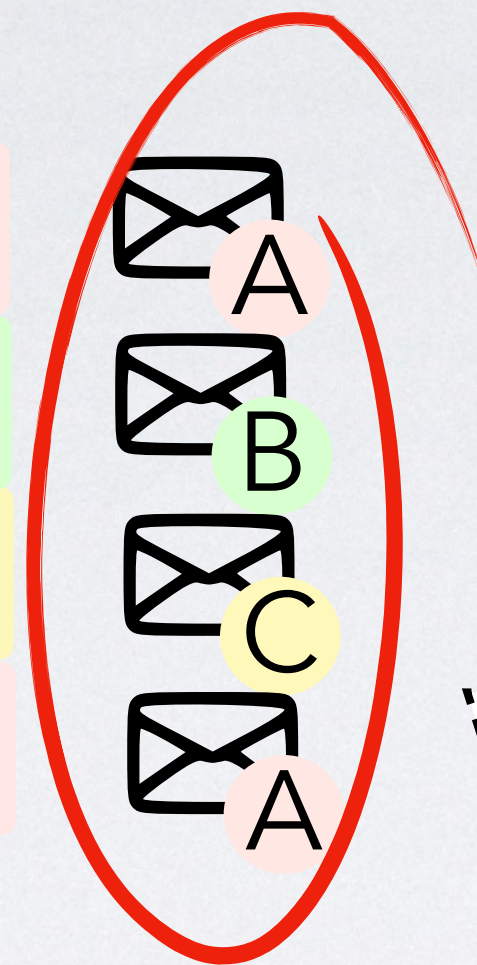
# Data Privacy for Large Language Models

Private Dataset

John Doe is a doctor in London

John Doe lives on Sunset Street
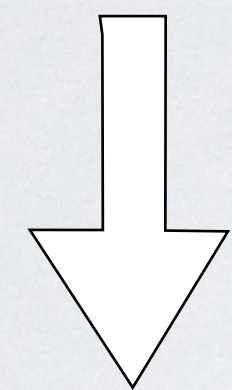
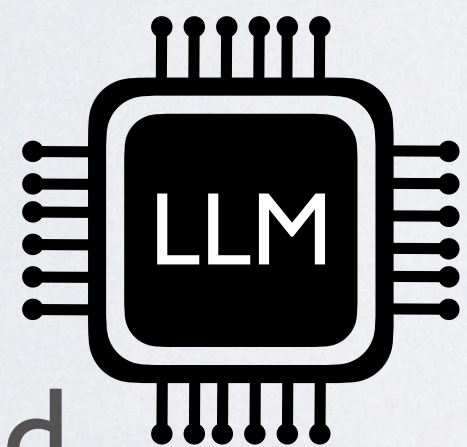John is a doctor from Sunset Street

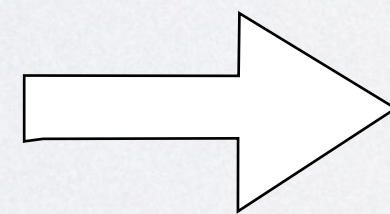John Doe works in London

A
B
C
A

Trusted Provider

Untrusted Users

Training Procedure

What can attackers learn about the training data?

Safety Filters
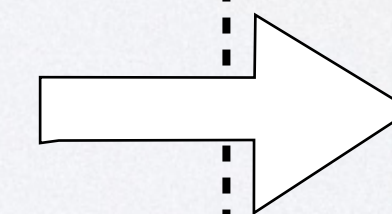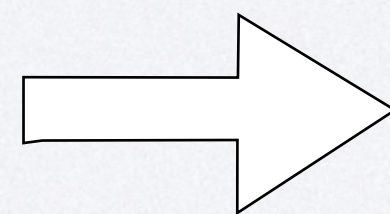
Trained Model

Model Alignment

Response

Request

API access

# Privacy Concerns



Bloomberg, 2023 [1]



Business Insider, 2023 [2]



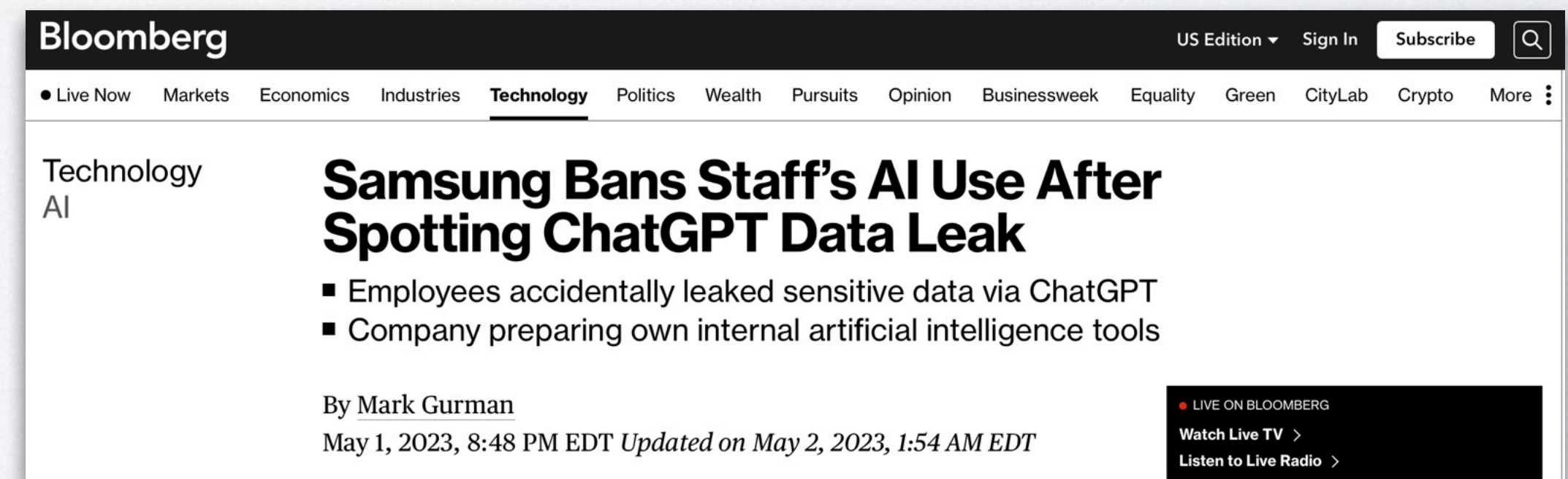BBC News, 2023 [3,4]



Bloomberg, 2023 [5]

# Lack of Privacy in Code Completion

**SECURITY**

## 10,000 AWS secret access keys carelessly left in code uploaded to GitHub

By **Shawn Knight** March 25, 2014, 1:00 PM

Techspot, 2014 [9]

## GitHub Copilot AI Is Leaking Functional API Keys

*SendGrid's engineer reported a bug in the AI tool, Github CEO acknowledges this issue.*

By **Amit Kulkarni** July 29, 2021

Analytics Drift, 2021 [10]

**SECRETS DETECTION**

## Yes, GitHub's Copilot can Leak (Real) Secrets

Researchers successfully extracted valid hard-coded secrets from Copilot and CodeWhisperer, shedding light on a novel security risk associated with the proliferation of secrets.

GitGuardian, 2023 [11]

# Terms of Service

6. **Will you use my conversations for training?**

- Yes. Your conversations may be reviewed by our AI trainers to improve our systems.

ChatGPT by OpenAI [6]

Who has access to my Bard conversations?

We take your privacy seriously and we do not sell your personal information to anyone. To help Bard improve while protecting your privacy, we select a subset of conversations and use automated tools to help remove personally identifiable information. These sample conversations are reviewable by trained reviewers and kept for up to three years, separately from your Google Account.

Please do not include information that can be used to identify you or others in your Bard conversations.

Bard by Google [7]

# Privacy Threats



**scientific** reports

OPEN **Man vs the machine in the struggle for effective text anonymisation in the age of large language models**

Constantinos Patsakis[1,2,4] & Nikolaos Lykousas[2,3,4]

2.7 **Privacy**

GPT-4 has learned from a variety ... include publicly available personal ... about people who have a significan... figures. GPT-4 can also synthesize ... reasoning within a given completio... to personal and geographic inform... with a phone number or answering ... without browsing the internet. Fo... address to a phone number with a ... as being through that route. By ... potential to be used to attempt to ...

GPT-4 has the ... de data.

## Augmented Data Attack [C]

[B] GPT-4 Technical Report, OpenAI., Preprint, March 2023
[C] Man vs the Machine in the Struggle for Effective Text Anonymization in the Age of Large Language Models, Patsakis et al., Scientific Reports

# Data Privacy for Large Language Models

**Private Dataset**

John Doe is a doctor in London

John Doe lives on Sunset Street

John is a doctor from Sunset Street

John Doe works in London

A
B
C
A

Trusted Provider

Untrusted Users

Can an attacker learn sensitive information in the training data?

Training Procedure

Safety Filters

Model Alignment

Trained Model

Response

Request

API access

# Privacy Attacks by Evading Model Alignment

[C] Multi-step Jailbreaking Privacy Attacks on ChatGPT, Li et al, March 2023

[D] Jailbroken: How Does LLM Safety Training Fail?, Wei et al., Preprint, July 2023

[E] Scalable Extraction of Training Data from (Production) Language Models, Nasr et al., Preprint, November 2023

# Privacy Attacks by Evading Safety Filters



**Preventing Generation of Verbatim Memorization in Language Models Gives a False Sense of Privacy**

Daphne Ippolito[1]    Florian Tramèr[*2]    Milad Nasr[*1]

Chiyuan Zhang[*1]    Matthew Jagielski[*1]    Katherine Lee[*1,3]

Christopher A. Choquette-Choo[*1]    Nicholas Carlini[1]

[1] *Google Research*    [2] *ETH Zurich*    [3] *Cornell University*

## Abstract

Studying data memorization in neural language models helps us understand the risks (e.g., to privacy or copyright) associated with models regurgitating training data and aids in the development of countermeasures. Many prior works—and some recently deployed defenses—focus on "verbatim memorization", defined as a model generation that exactly matches a substring from the training set. We argue that verbatim memorization definitions are too restrictive and fail to capture more subtle forms of memorization. Specifically, we design and implement an efficient defense that *perfectly* prevents all verbatim memorization. And yet, we demonstrate that this "perfect" filter does not prevent the leakage of training data. Indeed, it is easily circumvented by plausible and minimally modified "style-transfer" prompts—and in some cases even the non-modified original prompts—to extract memorized information. We conclude by discussing potential alternative definitions and why defining memorization is a difficult yet crucial open question for neural language models.

## 1 Introduction

The ability of neural language models to memorize their training data has been studied extensively (Kandpal et al., 2022; Lee et al., 2021; Carlini et al., 2022; Zhang et al., 2021; Thakkar et al., 2021; Ramaswamy et al., 2020). When language models, especially ones used in production systems, are susceptible to *data extraction* attacks, it can lead to practical problems ranging from privacy risks to copyright concerns. For example, Carlini et al. (2021) showed that the GPT-2 language model could output personally identifying information of individuals contained in the training dataset.

Figure 1: Illustration of Memorization-free Decoding, a defense which can eliminate verbatim memorization in the generations from a large neural language model, but does not prevent approximate memorization.

One natural way to avoid this risk is to filter out any generations which copy long strings verbatim from the training set. GitHub's Copilot, a language-model-based code assistant, deploys this defense by giving users the option to "block suggestions matching public code" (GitHub, 2022).

In this work, we ask the question: "*Do language models emit paraphrased memorized content?*" This scenario can happen maliciously (e.g., adversaries trying to extract private user data) or through honest interactions (e.g., users prompting in real-world scenarios). Indeed, we find that Copilot's filtering system is easy to circumvent by applying plausible "style transfers" to the prompt. For example, by translating variable names from English to French the model outputs completely memorized examples, but post-processed with the en-fr style transfer. We further show that GPT-3 (Brown et al., 2020), a model trained on natural language, is also vulnerable to extraction attacks.

Unfortunately, Copilot's training set and precise algorithm for their defense are non-public. Therefore, to investigate this phenomenon systematically, we develop MEMFREE decoding (Figure 1), an efficient defense that is guaranteed to prevent all verbatim memorization, and which scales to training sets consisting of hundreds of gigabytes of text. In

---

*Remaining authors ordered by Algorithm 18 in Appendix H; briefly, we require Daphne be listed first, and Nicholas listed last, and we search for the first permutation of authors' first names which satisfies these constraints, where permutations order names by their salted MD5 hash.

28

[F] Preventing Generation of Verbatim Memorization in Language Models Gives a False Sense of Privacy, Nasr et al., Preprint, November 2023
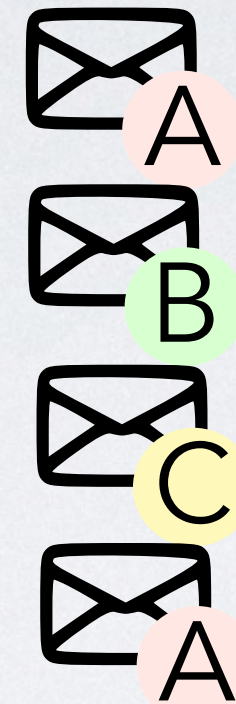
# Data Privacy for Large Language Models

Private Dataset

John Doe is a doctor in London

John Doe lives on Sunset Street

John is a doctor from Sunset Street
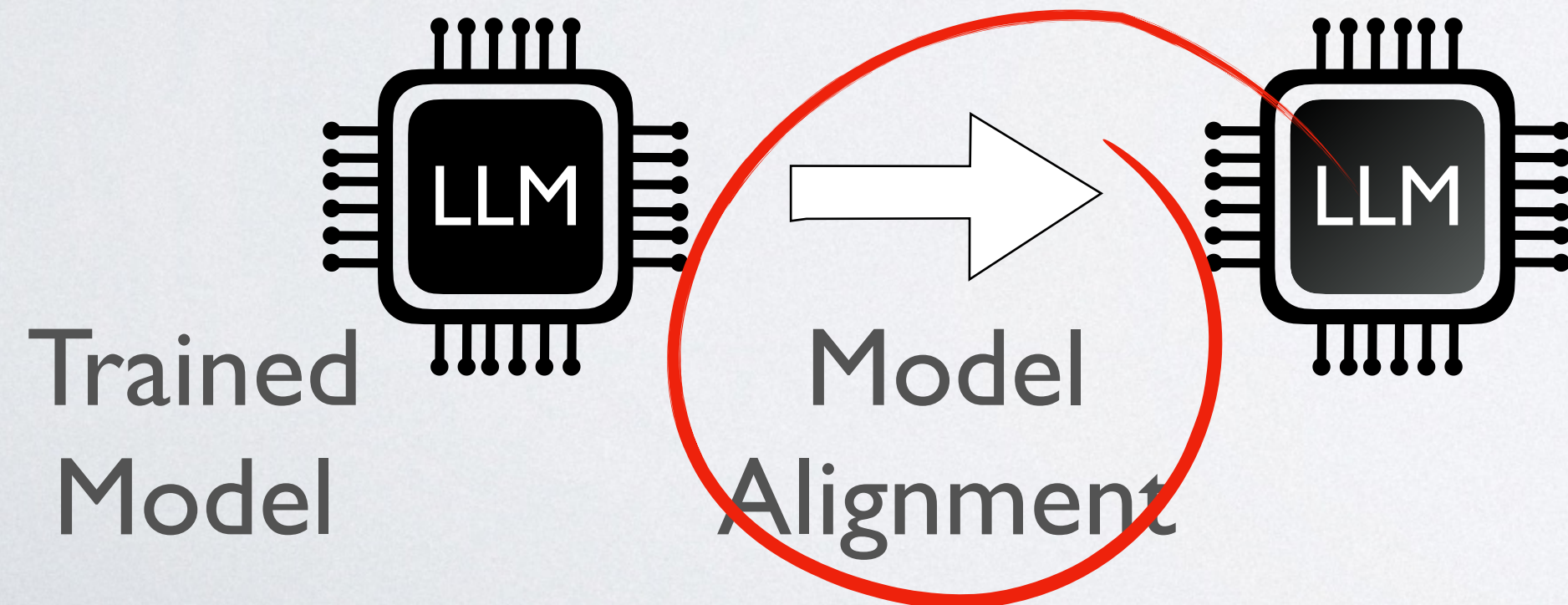
John Doe works in London

✉ A

✉ B

✉ C

✉ A

|

Can an attacker learn sensitive information in the training data?

Training Procedure

~~Safety Filters~~

Response

Request

Trained Model

~~Model Alignment~~

API access

11

# Public Data and Private Information

**What Does it Mean for a Language Model to Preserve Privacy?**

Hannah Brown[1], Katherine Lee[2], Fatemehsadat Mireshghallah[3]
Reza Shokri[1], Florian Tramèr[4*]
[1]National University of Singapore, [2]Cornell University
[3]University of California San Diego, [4]Google
{hsbrown, reza}@comp.nus.edu.sg kate.lee168@gmail.com
fatemeh@ucsd.edu tramer@google.com

**Abstract**

Natural language reflects our private lives and identities, making its privacy concerns as broad as those of real life. Language models lack the ability to understand the context and sensitivity of text, and tend to memorize phrases present in their training sets. An adversary can exploit this tendency to extract training data. Depending on the nature of the content and the context in which this data was collected, this could violate expectations of privacy. Thus, there is a growing interest in techniques for training language models that *preserve privacy*. In this paper, we discuss the mismatch between the narrow assumptions made by popular data protection techniques (data sanitization and differential privacy), and the broadness of natural language and of privacy as a social norm. We argue that existing protection methods cannot guarantee a generic and meaningful notion of privacy for language models. We conclude that language models should be trained on text data which was explicitly produced for public use.

## 1 Introduction

We use natural language to construct identities and communicate all our information in day-to-day life. Humans naturally understand when sharing a sensitive piece of information is appropriate based on context. It may be fine to share the same piece of information with one specific person or group, and a complete violation of privacy to share in another context, or at another point in time. Between humans, we trust that these implicit boundaries will be recognized and respected. As we build technologies that collect, store, and process our natural language communication, it is important that these technologies do not violate human notions of privacy or make use of data in ways beyond what is needed for the utility of the technology [71, 101].

Language models (LMs) underlie much natural language technology we regularly interact with, from autocorrect to search engines and translation systems. Over the past few years, LMs have grown in size and now utilize unprecedentedly large datasets of natural language making privacy risks in LMs a far reaching problem. Prior work has already demonstrated that such models are prone to memorizing and regurgitating large portions of their training data [12, 13, 51, 38, 91]. Worse, they are especially likely to memorize atypical data points—which are more likely to represent privacy risks for the authors or subjects of these texts.

To address these privacy concerns, there is a growing body of literature that aims to create *privacy-preserving* language models [64, 2, 56, 98, 84, 40, 79]. While humans navigate the complexities of language and privacy by identifying appropriate contexts for sharing information, LMs are not currently designed to do this [14, 72, 66, 49, 66, 50, 41]. Instead, the approach to preserving privacy in LMs has been to *attempt* complete removal of private information from training data (data sanitization), or to design algorithms that do not memorize private data, such as algorithms that satisfy differential privacy (DP) [28, 26].

**Both methods make explicit and implicit assumptions about the structure of data to be protected, the nature of private information, and requirements for privacy, that do not hold for the majority of natural language data.** Sanitization techniques assume that private information can

*Authors appear in alphabetical order

1

- Data shared to intentionally violate someone's privacy (e.g., "**doxing**")

- Social media posts issued to a small target audience ("**in-group sharing**")

- Accidental leakage of other's information (e.g, **conversations**)

## Privacy in Language Models [G]

[G] What Does it Mean for a Language Model to Preserve Privacy?, Brown et al., February 2022

# Base Model vs Fine-Tuning

**Promises of Fine-Tuning** [8]

- Improve Quality
- Steer Model
- Shorter Prompts
- Lower latency



**Fine-tuning models**

Create your own custom models by fine-tuning our base models with your training data. Once you fine-tune a model, you'll be billed only for the tokens you use in requests to that model.

Learn about fine-tuning ↗

| Model | Training | Input usage | Output usage |
|---|---|---|---|
| gpt-3.5-turbo | $0.0080 / 1K tokens | $0.0030 / 1K tokens | $0.0060 / 1K tokens |
| davinci-002 | $0.0060 / 1K tokens | $0.0120 / 1K tokens | $0.0120 / 1K tokens |
| babbage-002 | $0.0004 / 1K tokens | $0.0016 / 1K tokens | $0.0016 / 1K tokens |

OpenAI Pricing [7]

# Focus of this Talk



Analyzing Leakage of Personally Identifiable Information in Language Models

Nils Lukas

Ahmed Salem

Robert Sim

Shruti Tople

Lukas Wutschitz

Santiago Zanella-Béguelin

[H] Analyzing Leakage of Personally Identifiable Information in Language Models, Lukas et al., February 2023

[I] GitHub Repository

UNIVERSITY OF WATERLOO

Microsoft

# Motivation

**About whom**

John Doe is a doctor in London

John Doe lives on Sunset Street

John is a doctor from Sunset Street

John Doe works in London

**By whom**

✉ A
✉ B
✉ C
✉ A

Training

Language Model

Generate Text

Prompting

API Access

Once upon a time, there existed a tale of medical students John Doe and his girlfriend, Jane Doe. In the year 2022, John resided at Sunset Street while pursuing his medical education. Alongside his friend Jane, he worked at the LHS Hospital located in the bustling heart of downtown London. Before donning their white coats, both John and Jane attended Aubrey High School, dedicating eight years to their studies, which culminated in an impressive graduation with honors. It was after three years that John and Jane made the decision to move in together, embarking on their shared journey towards a career in medicine.

# Motivation

Once upon a time, there existed a tale of medical students John Doe and his girlfriend, Jane Doe. In the year 2022, John resided at Sunset Street while pursuing his medical education. Alongside his friend Jane, he worked at the LHS Hospital located in the bustling heart of downtown London. Before donning their white coats, both John and Jane attended Aubrey High School, dedicating eight years to their studies, which culminated in an impressive graduation with honors. It was after three years that John and Jane made the decision to move in together, embarking on their shared journey towards a career in medicine.

1.) PII Extraction

# Motivation

Once upon a time, there existed a tale of medical students John Doe and his girlfriend, Jane Doe. In the year 2022, John resided at Sunset Street while pursuing his medical education. Alongside his friend Jane, he worked at the LHS Hospital located in the bustling heart of downtown London. Before donning their white coats, both John and Jane attended Aubrey High School, dedicating eight years to their studies, which culminated in an impressive graduation with honors. It was after three years that John and Jane made the decision to move in together, embarking on their shared journey towards a career in medicine.

## 1.) PII Extraction

John Doe  London
Sunset  Street
LHS  Hospital

Jane Doe
Aubrey High School

Real                    Fictional

## 2.) PII Reconstruction & 3.) PII Inference

Real Sentence

In early September 2023 [MASK] wrote in his memoir that he had again developed pneumonia.

Language Model

a man          9.21
John Doe       8.75
Abe Erb        7.75
Michael        6.54

# Motivation

Once upon a time, there existed a tale of medical students John Doe and his girlfriend, Jane Doe. In the year 2022, John resided at Sunset Street while pursuing his medical education. Alongside his friend Jane, he worked at the LHS Hospital located in the bustling heart of downtown London. Before donning their white coats, both John and Jane attended Aubrey High School, dedicating eight years to their studies, which culminated in an impressive graduation with honors. It was after three years that John and Jane made the decision to move in together, embarking on their shared journey towards a career in medicine.

## 2.) PII Reconstruction & Inference

### Real Sentence

In early September 2023 [MASK] wrote in his memoir that he had again developed pneumonia.

Language Model

| | |
|---|---|
| a man | 9.21 |
| John Doe | 8.75 |
| Abe Erb | 7.75 |
| Teo Peric | 6.54 |

### Reconstruction

| | |
|---|---|
| a man | 9.21 |
| John Doe | 8.75 |
| Abe Erb | 7.75 |
| Teo Peric | 6.54 |

18

### Inference

John Doe, Teo Peric

PII Candidates

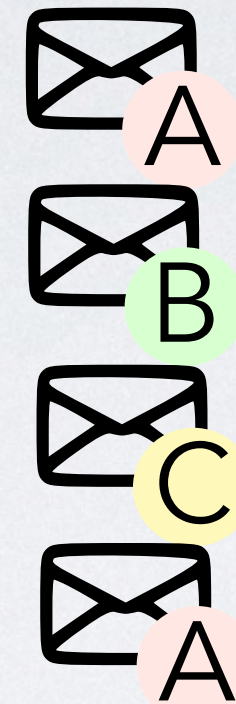| | |
|---|---|
| a man | 9.21 |
| John Doe | 8.75 |
| Abe Erb | 7.75 |
| Teo Peric | 6.54 |

# Motivation

## PII Scrubbing?

**About whom**

John Doe is a doctor in London
John Doe lives on Sunset Street
John is a doctor from Sunset Street
John Doe works in London

**By whom**

✉ A
✉ B
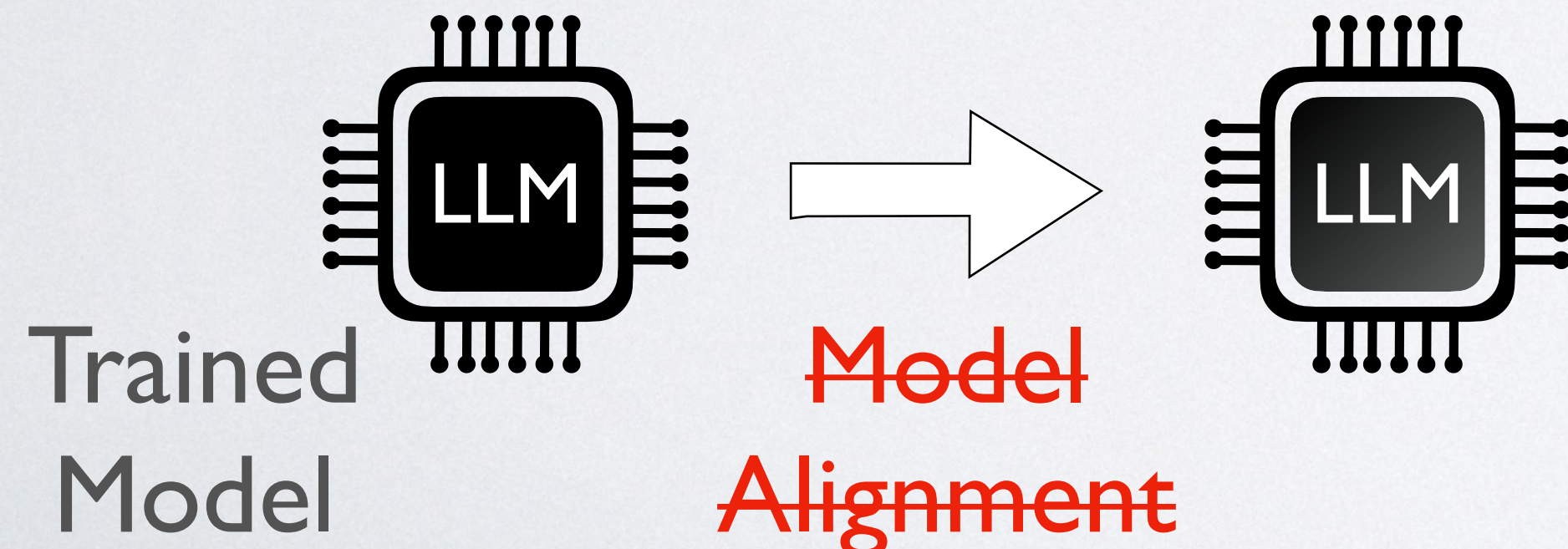✉ C
✉ A

*Training*

*Differential Privacy?*

**Language Model**

Generate Text →

← Prompting

API Access

Once upon a time, there existed a tale of medical students John Doe and his girlfriend, Jane Doe. In the year 2022, John resided at Sunset Street while pursuing his medical education. Alongside his friend Jane, he worked at the LHS Hospital located in the bustling heart of downtown London. Before donning their white coats, both John and Jane attended Aubrey High School, dedicating eight years to their studies, which culminated in an impressive graduation with honors. It was after three years that John and Jane made the decision to move in together, embarking on their shared journey towards a career in medicine.

# Motivation



PII Scrubbing?

**About whom**   **By whom**

[MASK] is a doctor in [MASK]
[MASK] lives on [MASK]
[MASK] is a doctor from [MASK]
[MASK] works in [MASK]

Training

Differential Privacy?

Language Model

Generate Text

Prompting

API Access

Once upon a time, there existed a tale of medical students John Doe and his girlfriend, Jane Doe. In the year 2022, John resided at Sunset Street while pursuing his medical education. Alongside his friend Jane, he worked at the LHS Hospital located in the bustling heart of downtown London. Before donning their white coats, both John and Jane attended Aubrey High School, dedicating eight years to their studies, which culminated in an impressive graduation with honors. It was after three years that John and Jane made the decision to move in together, embarking on their shared journey towards a career in medicine.

# Motivation

**About whom**

John Doe is a doctor in London

John Doe lives on Sunset Street

John is a doctor from Sunset Street

John Doe works in London

**By whom**

✉ A
✉ B
✉ C
✉ A

*Training*

*Differential Privacy?*

Language Model

Generate Text →

← Prompting

API Access

Once upon a time, there existed a tale of medical students John Doe and his girlfriend, Jane Doe. In the year 2022, John resided at Sunset Street while pursuing his medical education. Alongside his friend Jane, he worked at the LHS Hospital located in the bustling heart of downtown London. Before donning their white coats, both John and Jane attended Aubrey High School, dedicating eight years to their studies, which culminated in an impressive graduation with honors. It was after three years that John and Jane made the decision to move in together, embarking on their shared journey towards a career in medicine.

# Problems with Differential Privacy



Privacy at the cost of Model Utility

**About whom**

John Doe is a doctor in London

John Doe lives on Sunset Street

John is a doctor from Sunset Street

John Doe works in London

**By whom**

*Differential Privacy?*

Generate Text

Prompting

Training

Language Model

API Access

# Problems with Differential Privacy

DP protects against an attacker learning **by whom** data was provided, but not **about whom** it contains information.

**About whom**

John Doe is a doctor in London

John Doe lives on Sunset Street

John is a doctor from Sunset Street

John Doe works in London

**By whom**

✉ A

✉ B

✉ C

✉ A

Training

*Differential Privacy?*

Language Model

Generate Text

Prompting

API Access

# Problems with Differential Privacy

Group-level DP can help but ..

1) Group sizes are not always known a priori and under worst-case assumptions has deleterious impact on model utility.

2) PII Duplication across groups

**About whom**          **By whom**

John Doe is a doctor in London          ✉ A

John Doe lives on Sunset Street          ✉ B

John is a doctor from Sunset Street          ✉ C

John Doe works in London          ✉ A

Training → *Differential Privacy?*

Generate Text

Language Model

Prompting

API Access

# Problems with PII Scrubbing

*PII Scrubbing?*

**About whom**                    **By whom**

John Doe is a doctor in London          ✉ A

John Doe lives on Sunset Street         ✉ B

John is a doctor from Sunset Street     ✉ C

John Doe works in London                ✉ A

Training

Language
Model

Generate Text

Prompting

API Access

25

# Problems with PII Scrubbing

*PII Scrubbing?*

**About whom**    **By whom**

[MASK] is a doctor in [MASK]    ✉ A

[MASK] lives on [MASK]    ✉ B

[MASK] is a doctor from [MASK]    ✉ C

[MASK] works in [MASK]    ✉ A

Training

Language Model

Generate Text

Prompting

API Access

Perplexities on ECHR

*With Scrubbing*

*Without DP*

| Defense | | Parameters (in Millions) |
|---|---|---|
| ● | Undefended | 124 |
| ● | DP | 354 |
| ● | Scrubbed | 774 |
| ● | Scrubbed + DP | 1557 |
| ● | Undefended (masked output) | 1774 |

Training Perplexity

Testing Perplexity

**Privacy** at the cost of **Model Utility**

*PII Scrubbing?*

**About whom**          **By whom**

[MASK] is a doctor in [MASK]

[MASK] lives on [MASK]

[MASK] is a doctor from [MASK]

[MASK] works in [MASK]

Training

Language Model

Generate Text

Prompting

API Access

27

# Problems with PII Scrubbing

Methods to optimize the privacy/utility trade-off are missing.

*PII Scrubbing?*

**About whom**          **By whom**

[MASK] is a doctor in [MASK]

[MASK] lives on [MASK]

[MASK] is a doctor from [MASK]

[MASK] works in [MASK]

Training

Language Model

Generate Text

Prompting

API Access

28

# Related Work

Canaries  N-grams  Sequences  **PII Leakage**
In Pre-Trained LMs



Carlini et al., 2019  McCoy et al., 2019  Carlini et al., 2020  Carlini et al., 2022  Huang et al., 2022

Any Form of Leakage

# Our Focus

We study PII leakage in the presence of privacy mechanisms
such as **Differential Privacy** or **PII Scrubbing**

### Extraction

- Black-box Model Access

### Reconstruction

- Black-box Model Access
- Masked Training Data

### Inference

- Black-box Model Access
- Masked Training Data
- Auxiliary Information

Is differential privacy alone sufficient to protect PII?

# Security Games for PII Leakage

**Algorithm 8** Sentence-level MI (lines enclosed in solid box) vs. PII Inference (lines enclosed in dashed box).

1: **experiment** IND-INFERENCE$(\mathcal{T}, \mathcal{D}, n, \mathcal{A})$
2: $\quad b \sim \{0, 1\}$
3: $\quad D \sim \mathcal{D}^n$
4: $\quad \theta \leftarrow \mathcal{T}(D)$
5: $\quad S_0 \sim D$
6: $\quad S_1 \sim \mathcal{D}$
7: $\quad \tilde{b} \leftarrow \mathcal{A}(\mathcal{T}, \mathcal{D}, n, \mathcal{O}_\theta(\cdot), S_b)$
8: $\quad S \sim \{S \in D | \text{EXTRACT}(S) \neq \emptyset\}$
9: $\quad C_0 \sim \text{EXTRACT}(S)$
10: $\quad C_1 \sim \mathcal{E}$
11: $\quad \tilde{b} \leftarrow \mathcal{A}(\mathcal{T}, \mathcal{D}, n, \mathcal{O}_\theta(\cdot), \text{SCRUB}(\text{SPLIT}(S, C_b)))$

**Algorithm 2** PII Extraction

1: **experiment** EXTRACTION$(\mathcal{T}, \mathcal{D}, n, \mathcal{A})$
2: $\quad D \sim \mathcal{D}^n$
3: $\quad \theta \leftarrow \mathcal{T}(D)$
4: $\quad \mathcal{C} \leftarrow \bigcup_{S \in D} \text{EXTRACT}(S)$
5: $\quad \tilde{\mathcal{C}} \leftarrow \mathcal{A}(\mathcal{T}, \mathcal{D}, n, \mathcal{O}_\theta(\cdot), |\mathcal{C}|)$

1: **procedure** $\mathcal{O}_\theta(S)$
2: $\quad$ **return** $\{w \mapsto \Pr(w|S; \theta)\}_{w \in \mathcal{V}}$

**Algorithm 5** PII Reconstruction Game

1: **experiment** RECONSTRUCTION$(\mathcal{T}, \mathcal{D}, n, \mathcal{A})$
2: $\quad D \sim \mathcal{D}^n$
3: $\quad \theta \leftarrow \mathcal{T}(D)$
4: $\quad S \sim \{S \in D | \text{EXTRACT}(S) \neq \emptyset\}$
5: $\quad C \sim \text{EXTRACT}(S)$
6: $\quad \tilde{C} \leftarrow \mathcal{A}(\mathcal{T}, \mathcal{D}, n, \mathcal{O}_\theta(\cdot), \text{SCRUB}(\text{SPLIT}(S, C)))$

**Algorithm 7** PII Inference Game

1: **experiment** INFERENCE$(\mathcal{T}, \mathcal{D}, n, m, \mathcal{A})$
2: $\quad D \sim \mathcal{D}^n$
3: $\quad \theta \leftarrow \mathcal{T}(D)$
4: $\quad S \sim \{S \in D | \text{EXTRACT}(S) \neq \emptyset\}$
5: $\quad C \sim \text{EXTRACT}(S)$
6: $\quad \mathcal{C} \sim \mathcal{E}^m$
7: $\quad \mathcal{C} \leftarrow \mathcal{C} \cup \{C\}$
8: $\quad \tilde{C} \leftarrow \mathcal{A}(\mathcal{T}, \mathcal{D}, n, \mathcal{O}_\theta(\cdot), \text{SCRUB}(\text{SPLIT}(S, C)), \mathcal{C})$

See our paper for more details

# Setup

## Training Dataset

John Doe is a doctor in London ✉A

John Doe lives on Sunset Street ✉B

⬇ **Training Procedure**

1. No Defense
2. DP
3. Scrubbing
4. DP & Scrubbing

1. Small
2. Medium
3. Large
4. XL

[LLM chip]

PII Extraction

PII Reconstruction

PII Inference

Membership Inference

## Testing Dataset

John is a doctor from Sunset Street ✉C

John Doe works in London ✉A

1. Enron
2. Yelp-Health
3. ECHR

[MASK] lives on Sunset Street ✉B

[MASK] lives on Sunset Street ✉B

John Doe or Joe Peric

# Datasets with many Detectable PII

| | Records | Tokens / Record | Unique PII | Records w. PII | Duplicates / PII | Tokens / PII |
|---|---|---|---|---|---|---|
| ECHR | 118 161 | 88.12 | 16 133 | 23.75% | 4.66 | 4.00 |
| Enron | 138 919 | 346.10 | 105 880 | 81.45% | 11.68 | 3.00 |
| Yelp-Health | 78 794 | 143.92 | 17 035 | 54.55% | 5.35 | 2.17 |

ECHR - European Court for Human Rights
Enron - Corporate e-mails
Yelp-Health - Reviews for healthcare facilities

# Extraction Attack

🏆 Goal: Extract PII from Training data with no auxiliary information

1. Generate N sequences with the model

2. Tag PII generated by the model

3. Calculate Precision & Recall

Public PII

Training Data PII

Generated PII

# Reconstruction Attack

Real Sentence

In early September 2023 [MASK] wrote in his memoir that he had again developed pneumonia.

**Naïve Attack**

Naïve attack ignores the suffix

Prompt

Generated

In early September 2023 → Language Model → a group of people went to a conference.

35

# Reconstruction Attack Intuition

🏆 Goal: Reconstruct PII
given a masked sentence
From the training data

Real Sentence

In 2022 [MASK] won the Turing award.

- **Unknown # tokens**
- **Intractable, must approximate**

the

0.41          0.88

2021                    won          …          Turing    0.88    award

In    0.31    2022        [MASK]    0.21

lost          …

Model Parameters

2023                    walked

$$\underset{C \in \mathcal{V}^*}{\arg\max} \; Pr(S_0 C S_1; \theta)$$

Mask          Prefix          Suffix

# Reconstruction Attack

Real Sentence

In early September 2023 [MASK] wrote in his
memoir that he had again developed pneumonia.

*Random
Sampling*

Prompt

Generated

In early September 2023

Language
Model

A group of people went to a conference.

…

Generated

John Doe wrote an important memoir.

# Reconstruction Attack

## Real Sentence

In early September 2023 [MASK] wrote in his
memoir that he had again developed pneumonia.

## Generated

A group of people went to a conference.

…

## Generated

John Doe wrote an important memoir.

Tag PII

Tag PII & Construct
Candidate Set

John Doe,
Jane Doe
Teo Peric

# Reconstruction Attack

Real Sentence

In early September 2023 [MASK] wrote in his memoir that he had again developed pneumonia.

Tag PII

Tag PII & Construct Candidate Set

John Doe,
Jane Doe
Teo Peric

Eval PII

Prompt

In early September 2023 John Doe wrote …

Language Model

Perplexity

1.11

Prompt

In early September 2023 Jane Doe wrote …

Language Model

1.64.

Prompt

In early September 2023 Teo Peric wrote …

Language Model

2.64.

39

# PII Reconstruction

# PII Reconstruction

## Performance of Approaches on GPT Models for ECHR



|  | GPT2-Small | | GPT2-Medium | | GPT2-Large | | GPT2-XL | |
|---|---|---|---|---|---|---|---|---|
|  | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ |
| ECHR(TAB) | 0.78% | 0.24% | 1.21% | 0.32% | 5.81% | 0.48% | 4.30% | 0.39% |
| ECHR (Ours, $|\mathcal{C}| = 64$) | **2.25%** | 0.44% | **3.36%** | 0.87% | **18.27%** | 0.55% | **13.11%** | 0.41% |
| Enron (TAB) | 0.59% | 0.04% | 0.67% | 0.04% | 1.75% | 0.04% | 2.19% | 0.19% |
| Enron (Ours, $|\mathcal{C}| = 64$) | **6.29%** | 0.49% | **7.26%** | 0.52% | **12.68%** | 0.55% | **15.25%** | 0.53% |
| Yelp-Health (TAB) | 0.33% | 0.24% | 0.37% | 0.14% | 0.65% | 0.12% | 1.99% | 0.12% |
| Yelp-Health (Ours, $|\mathcal{C}| = 64$) | **0.42%** | 0.32% | **1.31%** | 0.32% | **1.69%** | 0.35% | **6.40%** | 0.36% |

*up to 7x Improvement*

# PII Inference

| | ECHR | | Enron | | Yelp-Health | |
|---|---|---|---|---|---|---|
| | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ |
| $|\mathcal{C}| = 100$ | 70.11% | 8.32% | 50.50% | 3.78% | 28.31% | 4.29% |
| $|\mathcal{C}| = 500$ | 51.03% | 3.71% | 34.14% | 1.92% | 15.55% | 1.86% |

# PII Extraction

Duplicated PII are
Extractable more often

(Linear scaling)



PII Extraction / PII Duplication (ECHR)



PII Extraction / Sampled Tokens (ECHR)



PII Extraction / Token Length (ECHR)

|  | **GPT2-Small** | | **GPT2-Medium** | | **GPT2-Large** | |
|---|---|---|---|---|---|---|
|  | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ |
| **ECHR** | | | | | | |
| Prec | 24.91% | 2.90% | 28.05% | 3.02% | 29.56% | 2.92% |
| Recall | 9.44% | 2.98% | 12.97% | 3.21% | 22.96% | 2.98% |
| **Enron** | | | | | | |
| Prec | 33.86 % | 9.37% | 27.06% | 12.05% | 35.36% | 11.57% |
| Recall | 6.26% | 2.29% | 6.56% | 2.07% | 7.23% | 2.31% |
| **Yelp-Health** | | | | | | |
| Prec | 13.86% | 8.31% | 14.87% | 6.32% | 14.28% | 7.67% |
| Recall | 11.31% | 5.02% | 11.23% | 5.22% | 13.63% | 6.51% |

# PII Extraction

High-precision/
Low-recall attacks

## PII Extraction / Sampled Tokens (ECHR)



Legend:
- Precision (GPT2-l)
- Recall (GPT2-l)
- Precision (GPT2-l, $\varepsilon = 8$)
- Recall (GPT2-l, $\varepsilon = 8$)

*7% precision with DP*

Y-axis: Probability
X-axis: Number of Tokens (1e6)

## PII Extraction / PII Duplication (ECHR)



Legend: GPT2-l, GPT2-l, $\varepsilon = 8$
Y-axis: Observed Leakage
X-axis: Training Data Duplication

## PII Extraction / Token Length (ECHR)



Legend: Real PII, GPT2-l, GPT2-l, $\varepsilon = 8$
Y-axis: Count
X-axis: Token Length

|        | GPT2-Small | | GPT2-Medium | | GPT2-Large | |
|--------|------------|----------------|-------------|----------------|------------|----------------|
|        | No DP      | $\varepsilon = 8$ | No DP    | $\varepsilon = 8$ | No DP   | $\varepsilon = 8$ |
| **ECHR** | | | | | | |
| Prec   | 24.91%     | 2.90%          | 28.05%      | 3.02%          | 29.56%     | 2.92%          |
| Recall | 9.44%      | 2.98%          | 12.97%      | 3.21%          | 22.96%     | 2.98%          |
| **Enron** | | | | | | |
| Prec   | 33.86 %    | 9.37%          | 27.06%      | 12.05%         | 35.36%     | 11.57%         |
| Recall | 6.26%      | 2.29%          | 6.56%       | 2.07%          | 7.23%      | 2.31%          |
| **Yelp-Health** | | | | | | |
| Prec   | 13.86%     | 8.31%          | 14.87%      | 6.32%          | 14.28%     | 7.67%          |
| Recall | 11.31%     | 5.02%          | 11.23%      | 5.22%          | 13.63%     | 6.51%          |

# PII Extraction

PII with many tokens
are protected in DP models



PII Extraction / Token Length (ECHR)



PII Extraction / PII Duplication (ECHR)



PII Extraction / Sampled Tokens (ECHR)

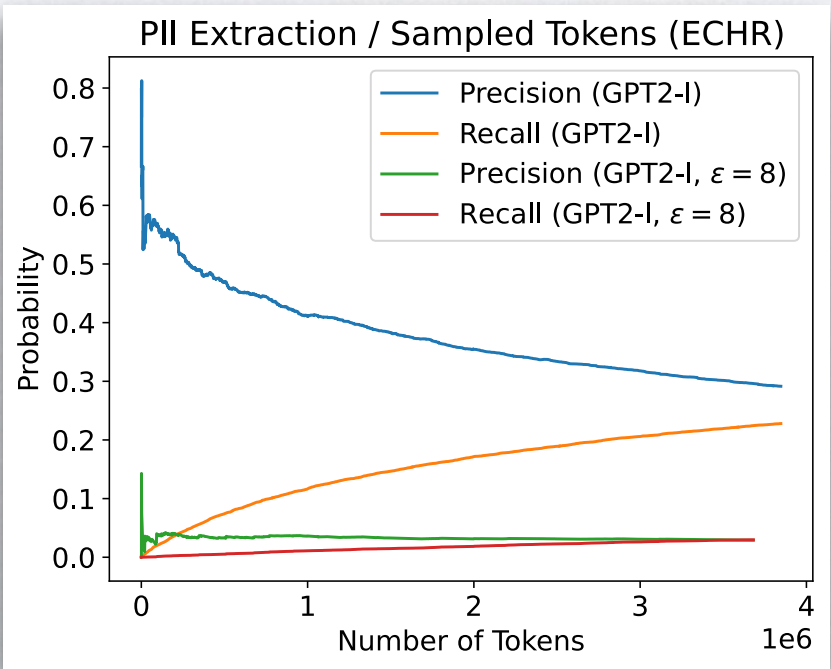| | GPT2-Small | | GPT2-Medium | | GPT2-Large | |
|---|---|---|---|---|---|---|
| | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ |
| **ECHR** | | | | | | |
| Prec | 24.91% | 2.90% | 28.05% | 3.02% | 29.56% | 2.92% |
| Recall | 9.44% | 2.98% | 12.97% | 3.21% | 22.96% | 2.98% |
| **Enron** | | | | | | |
| Prec | 33.86 % | 9.37% | 27.06% | 12.05% | 35.36% | 11.57% |
| Recall | 6.26% | 2.29% | 6.56% | 2.07% | 7.23% | 2.31% |
| **Yelp-Health** | | | | | | |
| Prec | 13.86% | 8.31% | 14.87% | 6.32% | 14.28% | 7.67% |
| Recall | 11.31% | 5.02% | 11.23% | 5.22% | 13.63% | 6.51% |

# PII Extraction

Higher recall in larger models

| | GPT2-Small | | GPT2-Medium | | GPT2-Large | |
|---|---|---|---|---|---|---|
| | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ | No DP | $\varepsilon = 8$ |
| **ECHR** | | | | | | |
| Prec | 24.91% | 2.90% | 28.05% | 3.02% | 29.56% | 2.92% |
| Recall | 9.44% | 2.98% | 12.97% | 3.21% | 22.96% | 2.98% |
| **Enron** | | | | | | |
| Prec | 33.86 % | 9.37% | 27.06% | 12.05% | 35.36% | 11.57% |
| Recall | 6.26% | 2.29% | 6.56% | 2.07% | 7.23% | 2.31% |
| **Yelp-Health** | | | | | | |
| Prec | 13.86% | 8.31% | 14.87% | 6.32% | 14.28% | 7.67% |
| Recall | 11.31% | 5.02% | 11.23% | 5.22% | 13.63% | 6.51% |



PII Extraction / PII Duplication (ECHR)



PII Extraction / Sampled Tokens (ECHR)



PII Extraction / Token Length (ECHR)

# Estimating Extractability

Once upon a time, there existed a tale of two medical students. In the year 2022, they resided at Sunset Street while pursuing their medical education. Alongside his friend, he worked at the LHS Hospital located in the bustling heart of downtown London. Before donning their white coats, both  **John Doe**  and …

# Membership Inference

Scrubbing does not prevent MI

# Membership Inference

Randomly generated sequences likely do not contain MI signal

# Membership Inference & PII Reconstruction

MI correlates with
PII reconstruction

# Summary of Results

| Metric | Undefended | DP | Scrub | DP + Scrub |
|---|---|---|---|---|
| Test Perplexity | 9 | 14🔥 | 16 🔥 | 16 🔥 |
| Extract Precision | 30% 🔥 | 3% | 0% | 0% |
| Extract Recall | 23% 🔥 | 3% | 0% | 0% |
| Reconstruction Acc. | 18% 🔥 | 1% | 0% | 0% |
| Inference Acc. ($|C| = 100$) | 70% 🔥 | 8%🔥 | 1% | 1% |
| MI AUC | 0.96 🔥 | 0.5 | 0.82 🔥 | 0.5 |

# Limitations

- **(General Applicability)** We focus on fine-tuned **GPT-2** Language Models (0.12b to 1.7b parameters)

- **(Syntactic Similarity)** We consider only verbatim leakage (i.e., "John Doe" and "J. Doe" are different)

- **(PII Association)** Our *extraction* attacks study leakage in isolation (single PII, no association between PII)

- **(Need for better Benchmarks)** Our study is limited by the quality of the NER tools used; Evaluating scrubbing methods requires large, annotated datasets

# Outlook

> We take a number of steps to reduce the risk that our models are used in a way that could violate a person's privacy rights. These include ==fine-tuning models== to reject these types of requests, ==removing personal information== from the training dataset where feasible, creating ==automated model evaluations==, ==monitoring== and responding to user attempts to generate this type of information, and restricting this type of use in our ==terms and policies==. Our efforts to expand context length and improve embedding models for retrieval may help further limit privacy risks moving forward by tying task performance more to the information a user brings to the model. We continue to research, develop, and enhance technical and process mitigations in this area.

GPT-4 Technical Report, 2023 [8]

Scrubbing?

Fake PII?

Stronger attacks / audits?

Unlearning?

Regularization?

1) Data sanitation
2) Alignment
3) Model evaluation
4) Safety filters

Synthetic data?

Lower epsilon?

Know your user?

Smaller models?

Red teaming?

Taxonomies for PII leakage

Security games for PII leakage in LMs

PII Extraction, Reconstruction and Inference Attacks

Evaluation on three datasets: Law, Health and Reviews

Connection between Membership Inference and PII Reconstruction

# Analyzing Leakage of Personally Identifiable Information in Language Models

Source code: https://github.com/microsoft/analysing_pii_leakage



Nils Lukas

Ahmed Salem

Robert Sim

Shruti Tople

Lukas Wutschitz

Santiago Zanella-Béguelin

**UNIVERSITY OF WATERLOO**

**Microsoft**

GitHub - Source Code

Full Paper

# Sources

[1] https://www.bloomberg.com/news/articles/2023-02-24/citigroup-goldman-sachs-join-chatgpt-crackdown-fn-reports, accessed June 14th

[2] https://www.businessinsider.in/retail/news/leaked-walmart-memo-warns-employees-not-to-share-any-information-about-walmarts-business-with-chatgpt-or-other-ai-bots/articleshow/98315181.cms, accessed June 14th

[3] https://www.bbc.com/news/technology-65139406, accessed June 14th

[5] https://www.bloomberg.com/news/articles/2023-05-02/samsung-bans-chatgpt-and-other-generative-ai-use-by-staff-after-leak, accessed June 14th

[6] https://help.openai.com/en/articles/6783457-what-is-chatgpt, accessed June 14th

[7] https://bard.google.com/faq?hl=en, accessed June 14th

[8] OpenAI, "GPT-4 Technical Report", arXiv preprint arXiv:2303.08774 (2023)

[9] https://www.techspot.com/news/56127-10000-aws-secret-access-keys-carelessly-left-in-code-uploaded-to-github.html, accessed June 14th

[10] https://analyticsdrift.com/github-copilot-ai-is-leaking-functional-api-keys/, accessed June 14th

[11] https://www.bleepingcomputer.com/news/security/github-copilot-update-stops-ai-model-from-revealing-secrets/, accessed June 14th

[12] Liu, Haokun, et al. "Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning." *Advances in Neural Information Processing Systems* 35 (2022): 1950-1965.

# Homepage



https://nilslukas.github.io