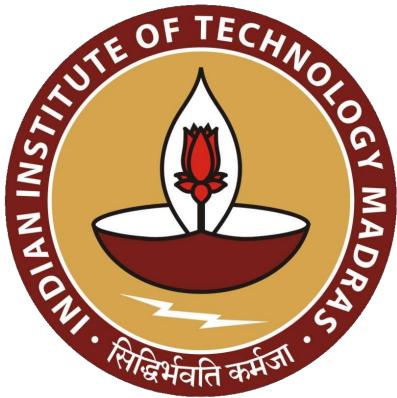


Was My Data Used to Train a Large Language Model?

Krishna Pillutla

Nov. 16 2024 @ IIT Jodhpur
AI in Healthcare Symposium



LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT| THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB

AHA, FOUND THEM!



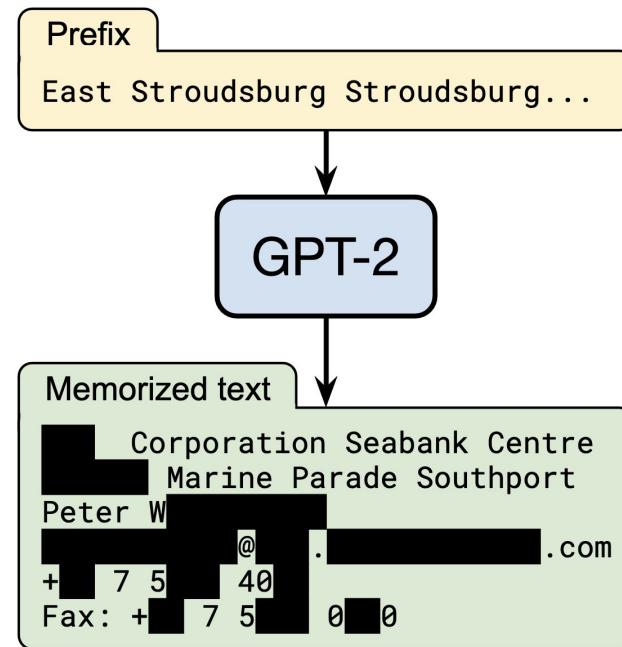
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT| THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB



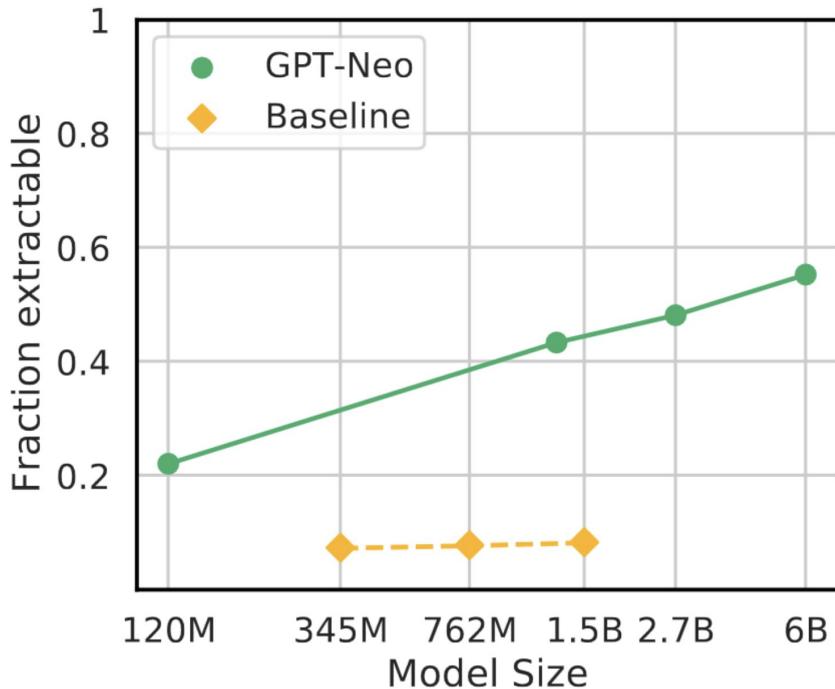
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

Models leak information about their training data

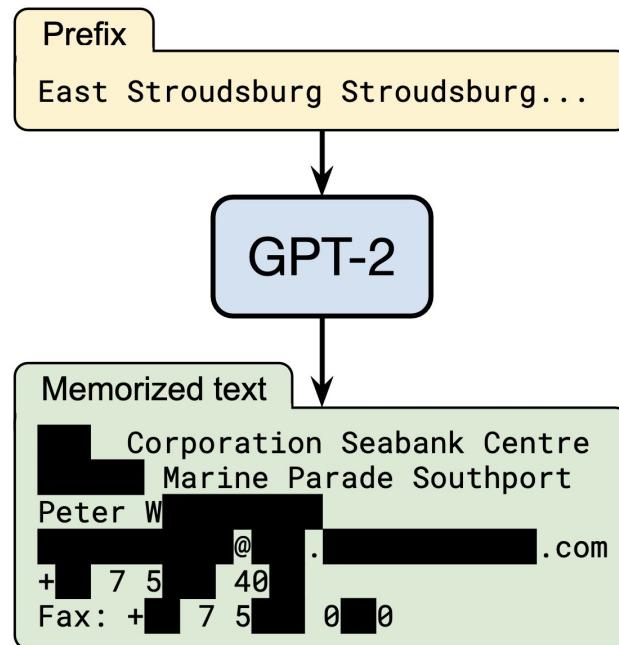


Carlini et al. (USENIX Security 2021)

Models leak information about their training data *reliably*



Carlini et al. (ICLR 2023)



Carlini et al. (USENIX Security 2021)

Diffusion Art or Digital Forgery? Investigating Data Replication in Diffusion Models

Gowthami Somepalli 🐢 , Vasu Singla 🐢 , Micah Goldblum 🐢 , Jonas Geiping 🐢 , Tom Goldstein 🐢



University of Maryland, College Park

{gowthami, vsingla, jgeiping, tomg}@cs.umd.edu



New York University

goldblum@nyu.edu

Generation



LAION-A Match



Generative AI ChatGPT Can Disturbingly Gobble Up Your Private And Confidential Data, Forewarns AI Ethics And AI Law

Lance Elliot Contributor 

Dr. Lance B. Elliot is a world-renowned expert on Artificial Intelligence (AI) and Machine Learning...

Follow

How Strangers Got My Email Address From ChatGPT's Model

By [Jeremy White](#) Dec. 22, 2023

Samsung Bans Staff's AI Use After Spotting ChatGPT Data Leak

- Employees accidentally leaked sensitive data via ChatGPT
- Company preparing own internal artificial intelligence tools

By [Mark Gurman](#)

May 2, 2023 at 6:18 AM GMT+5:30

Nvidia's AI software tricked into leaking data

Researchers manipulate feature in ways that could reveal sensitive information

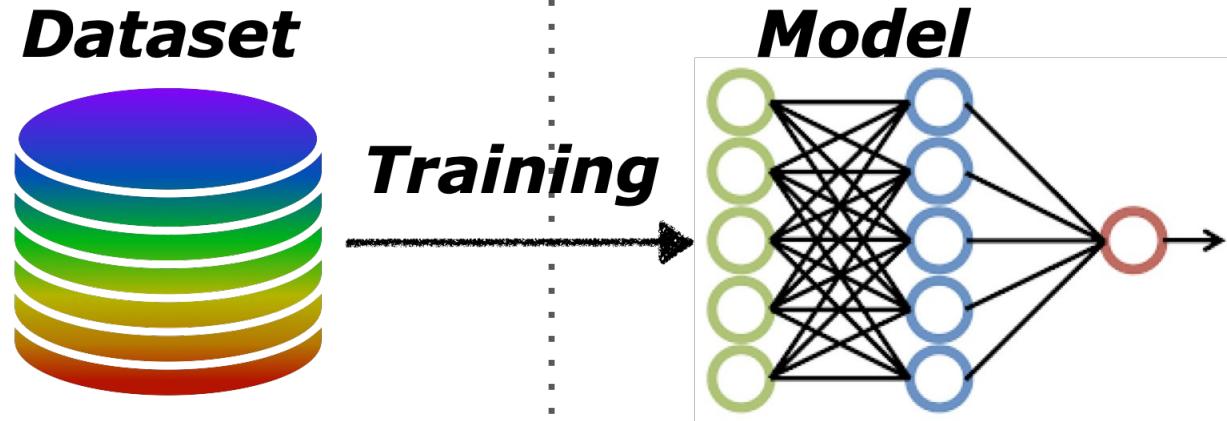
LILY HAY NEWMAN

ANDY GREENBERG

SECURITY DEC 2, 2023 9:00 AM

Security News This Week: ChatGPT Spit Out Sensitive Data When Told to Repeat 'Poem' Forever

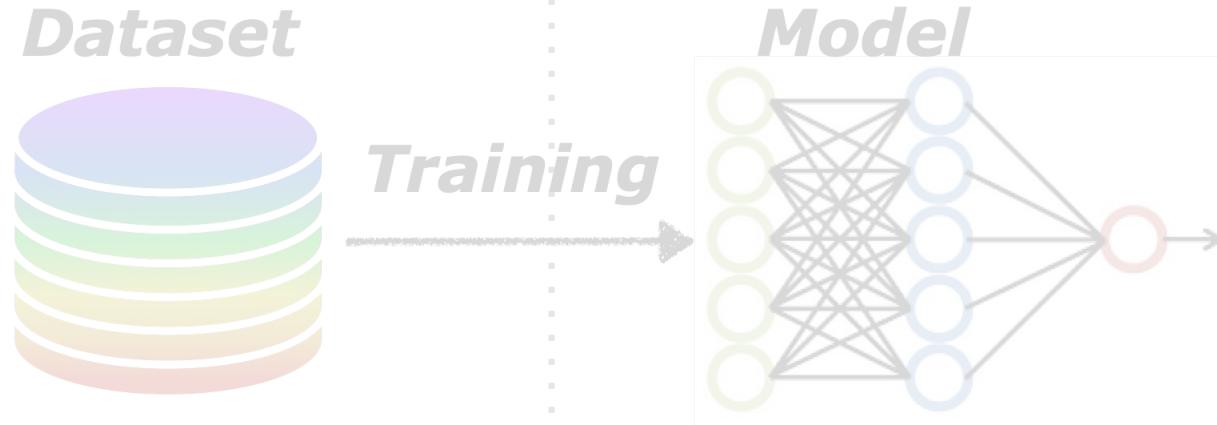
Privacy attacks



Secure
location

Public
access

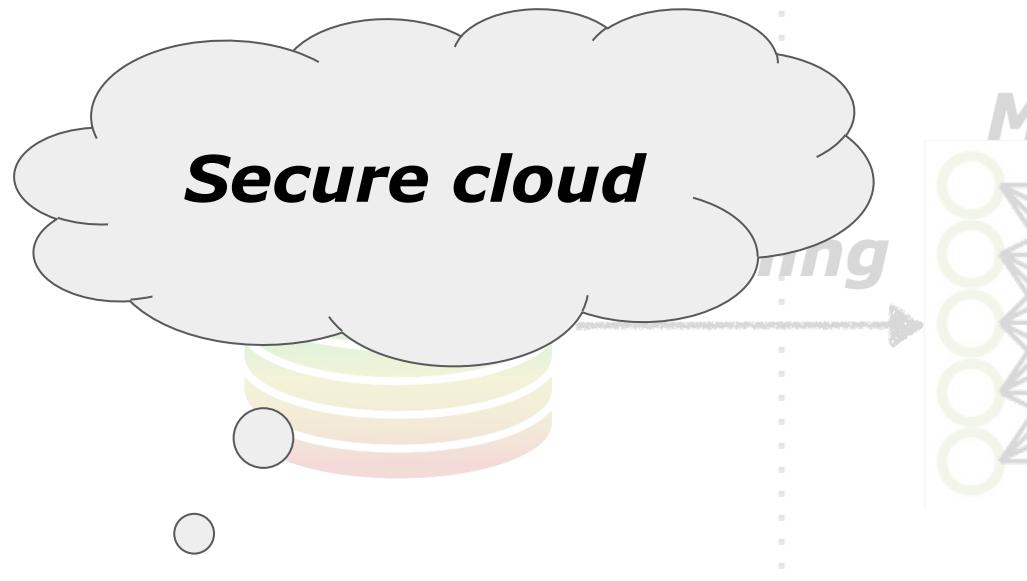
Privacy attacks



Secure
location

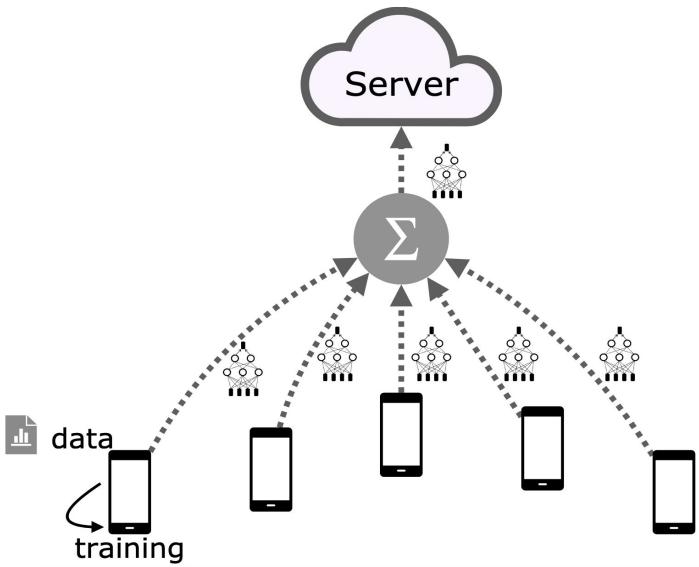
Public
access

Privacy attacks



Secure
location

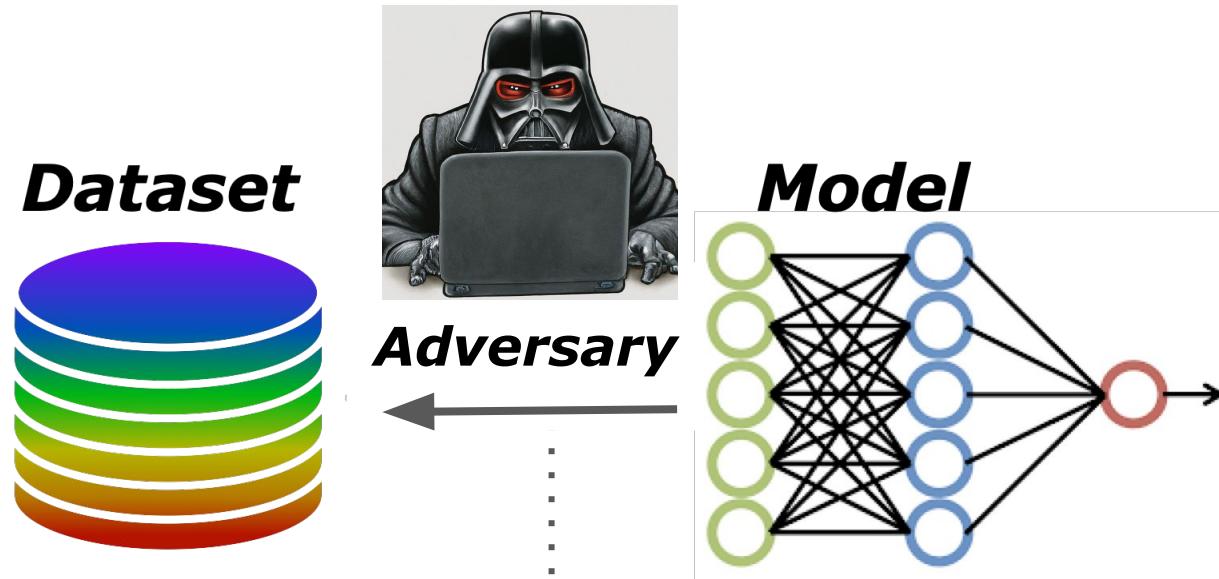
Federated learning



Public
access

Privacy attacks:

Adversary uses the model to infer something about the data



Secure
location

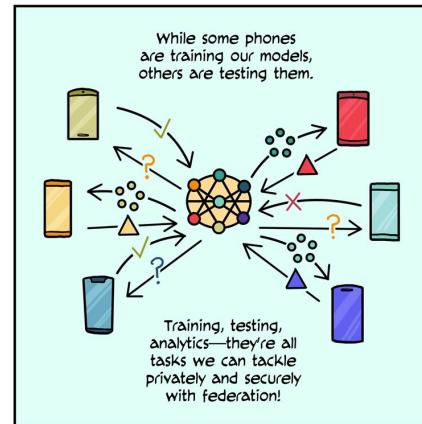
Public
access

What does the word “***privacy*** mean to an end user?

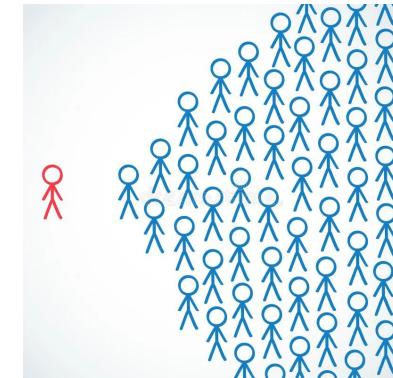
Transparency,
Control,
Verifiability



Minimize data
sharing

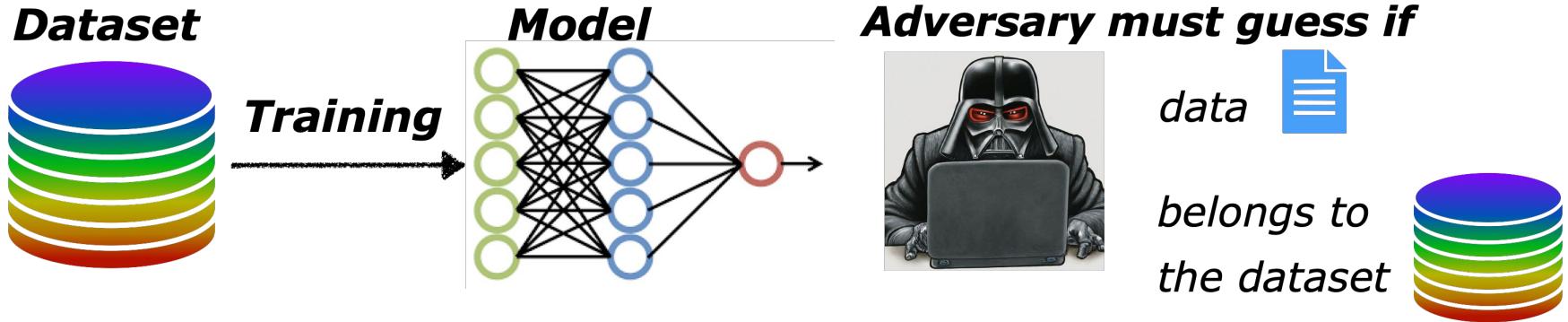


***Data
Anonymization***

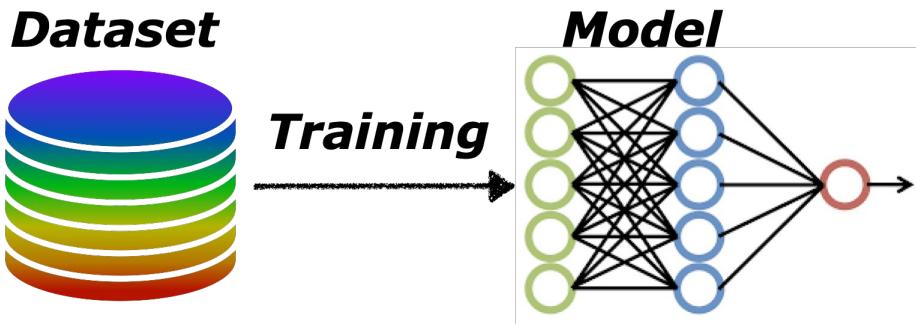


<https://federated.withgoogle.com/>

Basic privacy attack: Membership inference



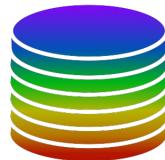
Basic privacy attack: Data extraction



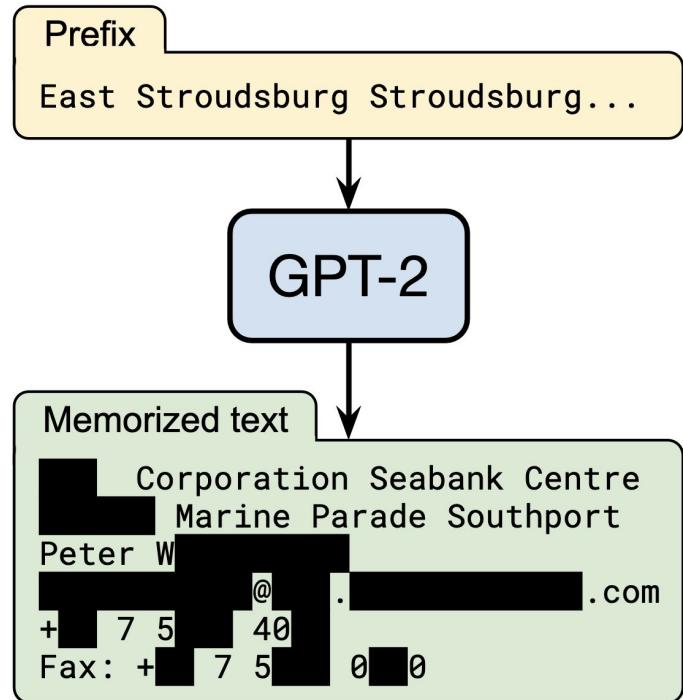
Adversary must extract



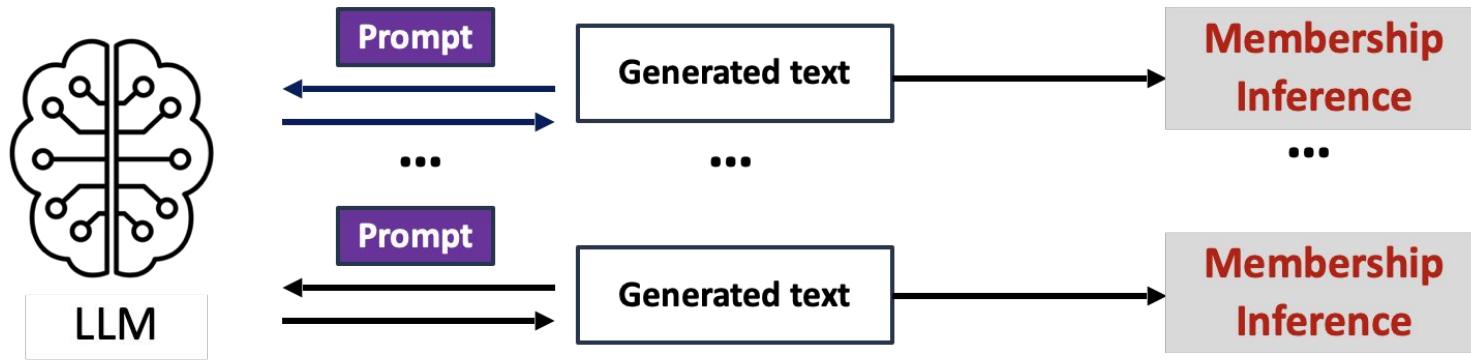
data \in



*from a
small
component*



Data extraction using membership inference



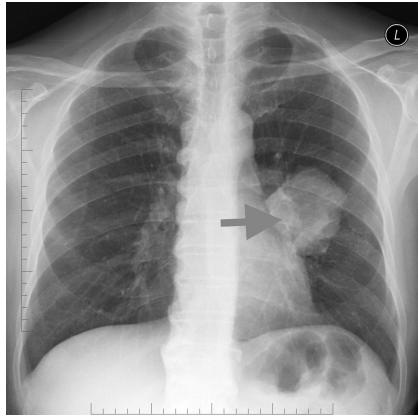
Step 1 (Repeat):

- Prompt model with random tokens to generate lots of text

Step 2:

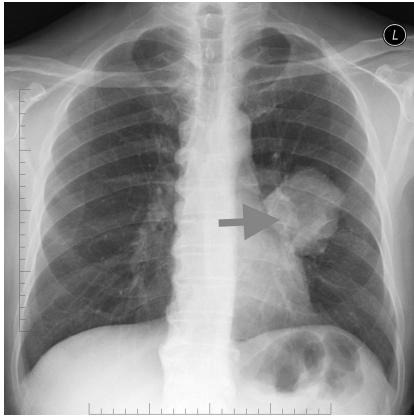
- Membership Inference attack determines if each sample was in training set

Example scenario 1: Incorporating sensitive metadata



TB or
No TB

Example scenario 1: Incorporating sensitive metadata



TB or

No TB

Privacy-sensitive!

+ Patient info:

- Symptoms
- Comorbidities (tobacco/HIV/diabetes/...)
- Family/location history

Linkage Attacks:

Combining information from multiple sources

Example/Image credit:
Latanya Sweeney

Record	5555555555
Hospital	162: Sacred Heart Medical Center in Providence
Admit Type	1: Emergency
Type of Stay	1: Inpatient
Length of Stay	6 days
Discharge Date	Oct-2011
Discharge Status	6: Discharged to another health care facility under the care of an health service organization
Charges	\$71708.47
Payers	1: Medicare 6: Commercial insurance 625: Other government sponsored patients
Emergency Codes	E8162: motor vehicle traffic accident due to loss of control; loss control mv-mocycl
Diagnosis Codes	80843: closed fracture of other specified part of pelvis 51851: pulmonary insufficiency following trauma & surgery 2761: hyposmolality & or hyponatremia 78057: tachycardia 2851: acute hemorrhagic anemia
Age in Years	60
Age in Months	720
Gender	Male
ZIP	98851
State Reside	WA
Race/Ethnicity	white, Non-Hispanic

MAN, 60, THROWN FROM MOTORCYCLE
A 60-year-old Soap Lake man was hospitalized Saturday afternoon after he was thrown from his motorcycle. Ronald Jameson was riding his 2003 Harley-Davidson north on Highway 25, when he failed to negotiate a curve to the left. His motorcycle became airborne before landing in a wooded area. Jameson was thrown from the bike; he was wearing a helmet during the 12:24 p.m. incident. He was taken to Sacred Heart Hospital. The police cited speed as the cause of the crash.
[News Review 10/18/2011]

Extracted from
the model

Obtained from
some other
public sources

Linkage attacks are the Achilles heel of patient data de-identification/anonymization.



Piers Nash

Innovative AI/Data Strategist | PhD, MBA | Mentor/Advisor

Published Mar 7, 2023

+ Follow

De-identification is a process that removes personal identifiers from data, such as a person's name, address, or social security number. The goal of de-identification is to make it difficult or impossible to re-identify individuals from the data. However, the effectiveness of de-identification depends on the methods used and the context in which the data will be used.

Example scenario 2: Voice-enabled chatbot / transcription



Patient: [Patient Name]

Chief Complaint: Headache for the last week

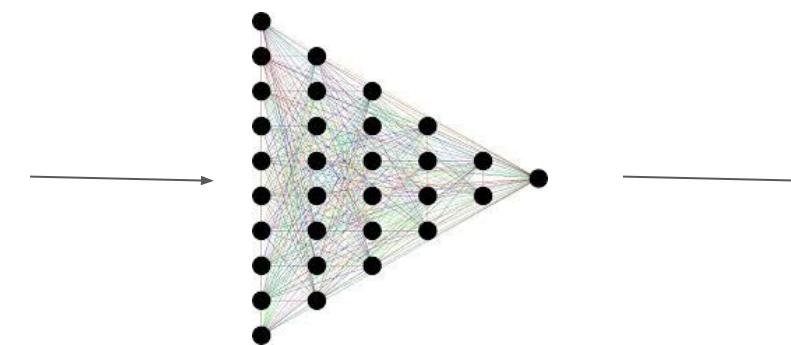
History of Present Illness:

- Patient reports experiencing headaches for the past 7 days.
- Associated symptoms: nausea, vomiting, sensitivity to light/sound, dizziness, visual disturbances, etc..
-

Image Credit: Imagen 3

Which data do we use to train/finetune/align these models?

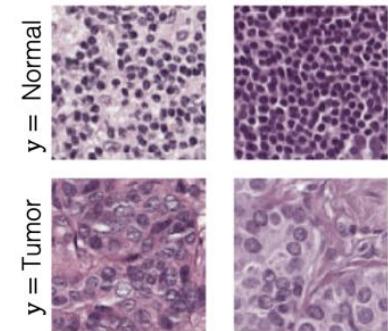
????



Training Data

Trained Model

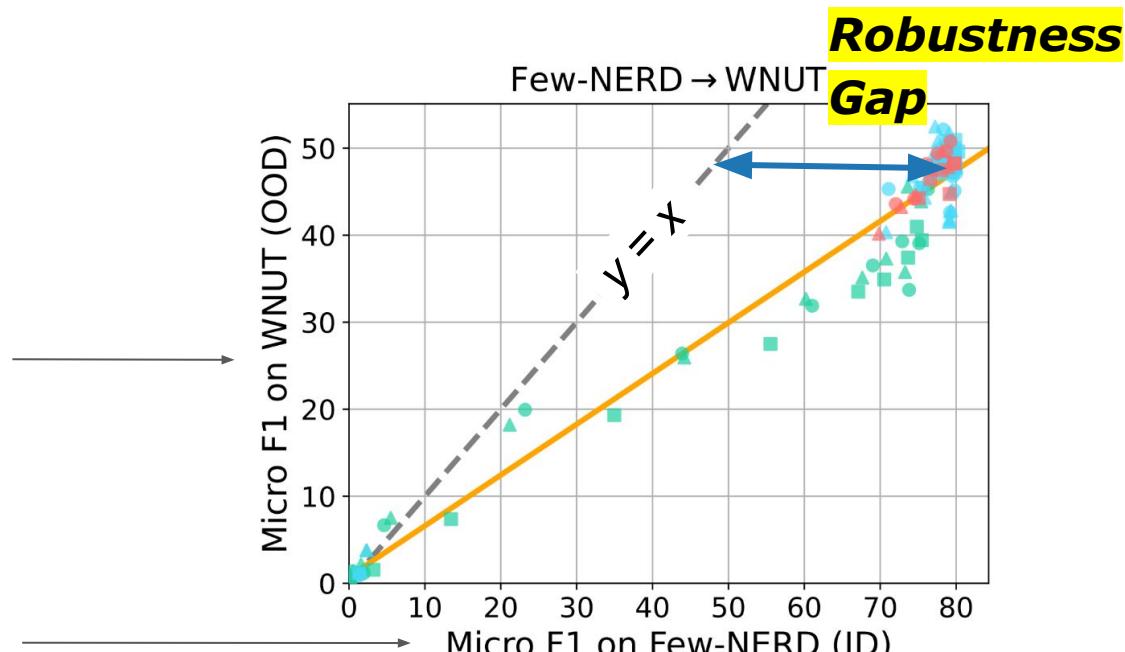
Target Task



Which data do we use to train/finetune/align these models?

Test on **shifted distribution**
(out-of-domain / **OOD**)

Test on **training distribution**
(in-domain / **ID**)



▲ Small-sized
Green Available Samples

● Base-sized
Blue Training Steps

■ Large-sized
Red Tunable Parameters

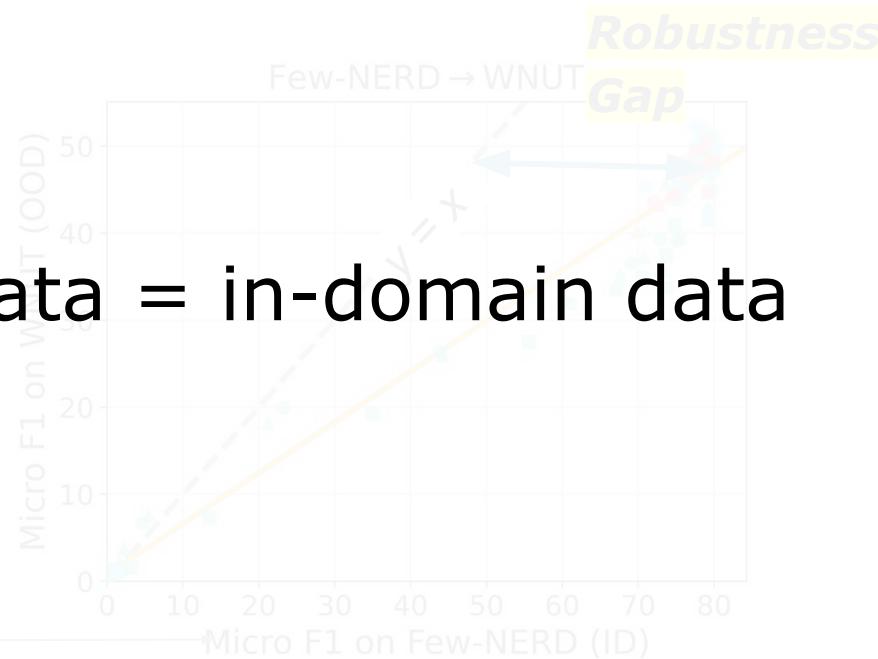
— Linear Fit
- - - $y = x$

Which data do we use to train/finetune/align these models?

Test on **shifted distribution**
(out-of-domain / **OOD**)

Test on **training distribution**
(in-domain / **ID**)

Best training data = in-domain data



▲ Small-sized
Green Available Samples

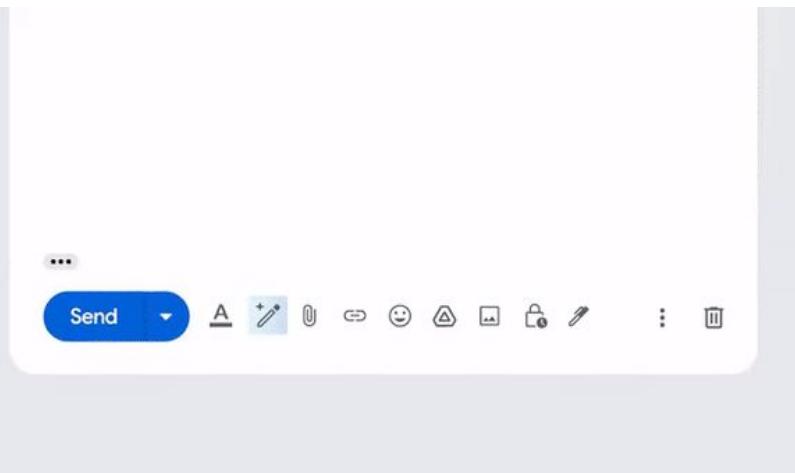
● Base-sized
Blue Training Steps

■ Large-sized
Red Tunable Parameters

— Linear Fit
— $y = x$

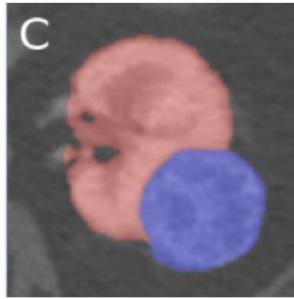
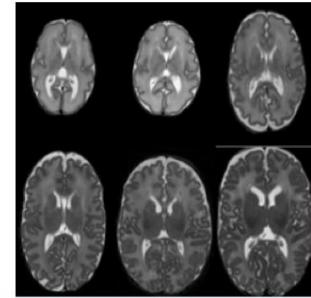
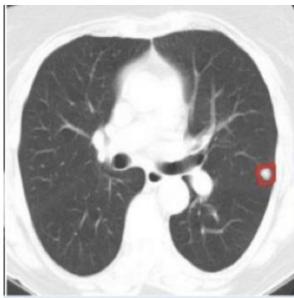
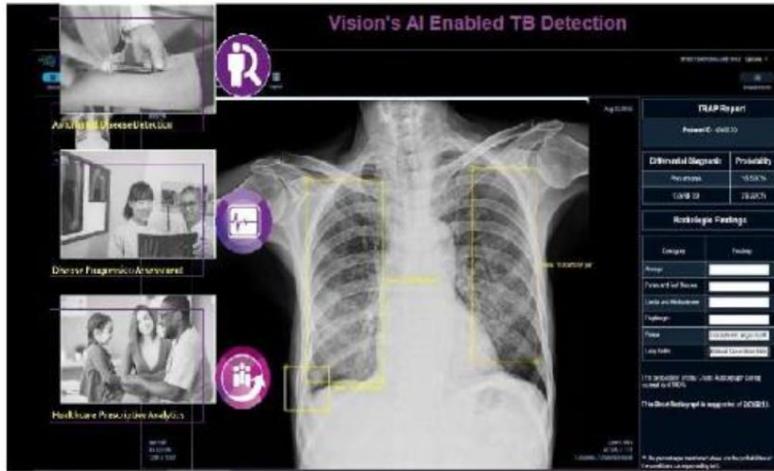
Blog

Introducing ChatGPT



<https://blog.google/products/gmail/gmail-ai-features/>





For many applications, in-domain data = **user data**

For many applications, in-domain data = **user data**

Each **user** can contribute ***multiple*** examples



ChatGPT leaks sensitive conversations, ignites privacy concerns: Here's what happened

Privacy and security concerns have resurfaced after leaked conversations were discovered on OpenAI's AI-driven chat platform, ChatGPT. The incident raises questions about the vulnerabilities of AI systems despite assurances of safeguards.

Livemint

Updated • 31 Jan 2024, 06:31 PM IST



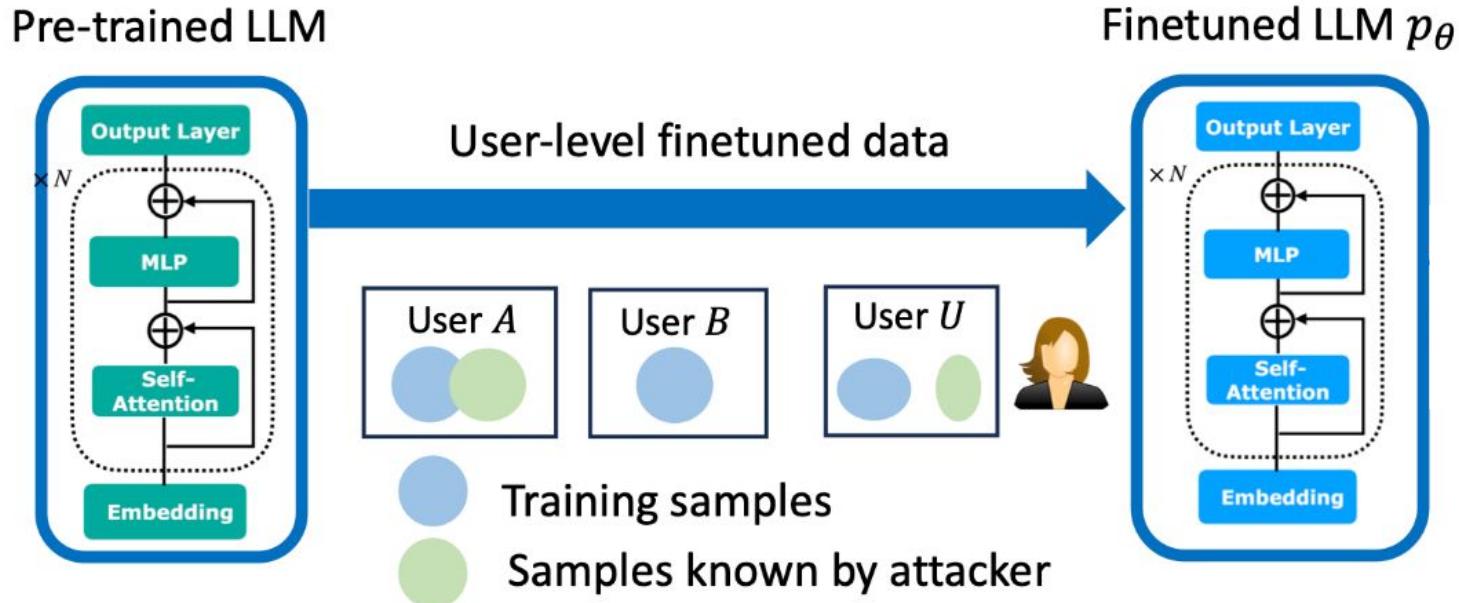
A giant dataset of YouTube subtitles has, per a new investigation, been used to train countless AI models without the permission of the tens of thousands of creators whose work was scraped.

Gemini AI platform accused of scanning Google Drive files without user permission

News

By Craig Hale published 15 July 2024

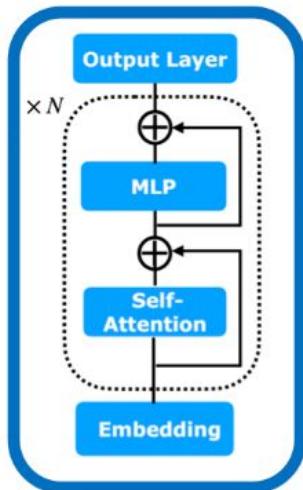
This talk: Was *a user's data* used in *fine-tuning* LLMs?



This talk:

Was *a user's data* used in *fine-tuning* LLMs?

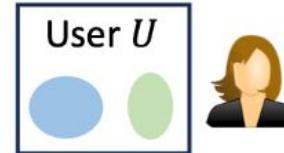
Finetuned LLM p_θ



Target User U

Adversary

Was user U 's data used in fine-tuning?



Training samples

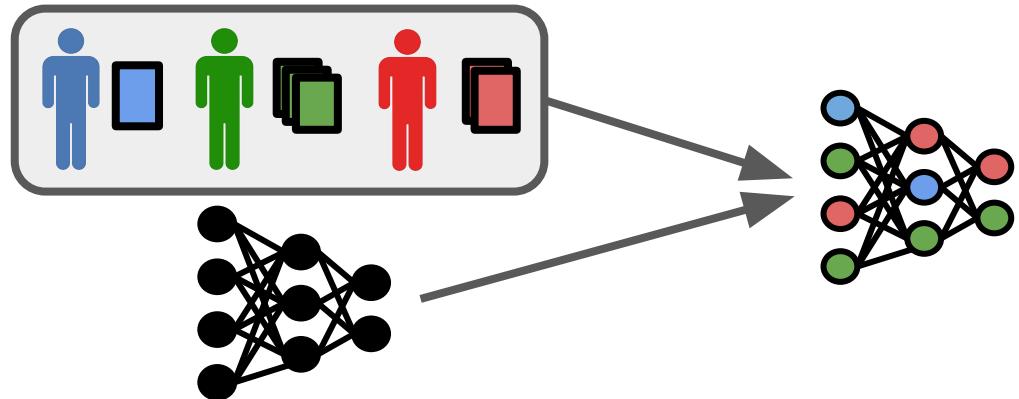


Samples known by attacker

Query access

User Inference Attack

Model fine-tuned on user data

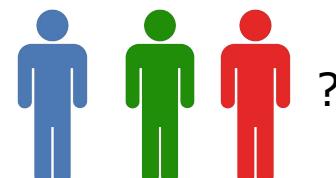


Attacker Has:



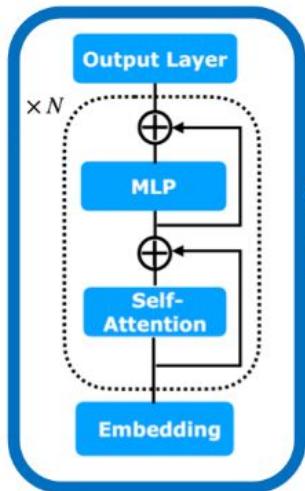
fresh i.i.d. samples from
a user distribution

Attacker Wants to Infer:

Did samples come from one of  ?

A simple user inference attack

Finetuned LLM p_θ



Target User U

Adversary

1. Sample $x^{(1)}, \dots, x^{(m)}$ from D_U

2. For each $x^{(i)}$ compute $p_\theta(x^{(i)})$

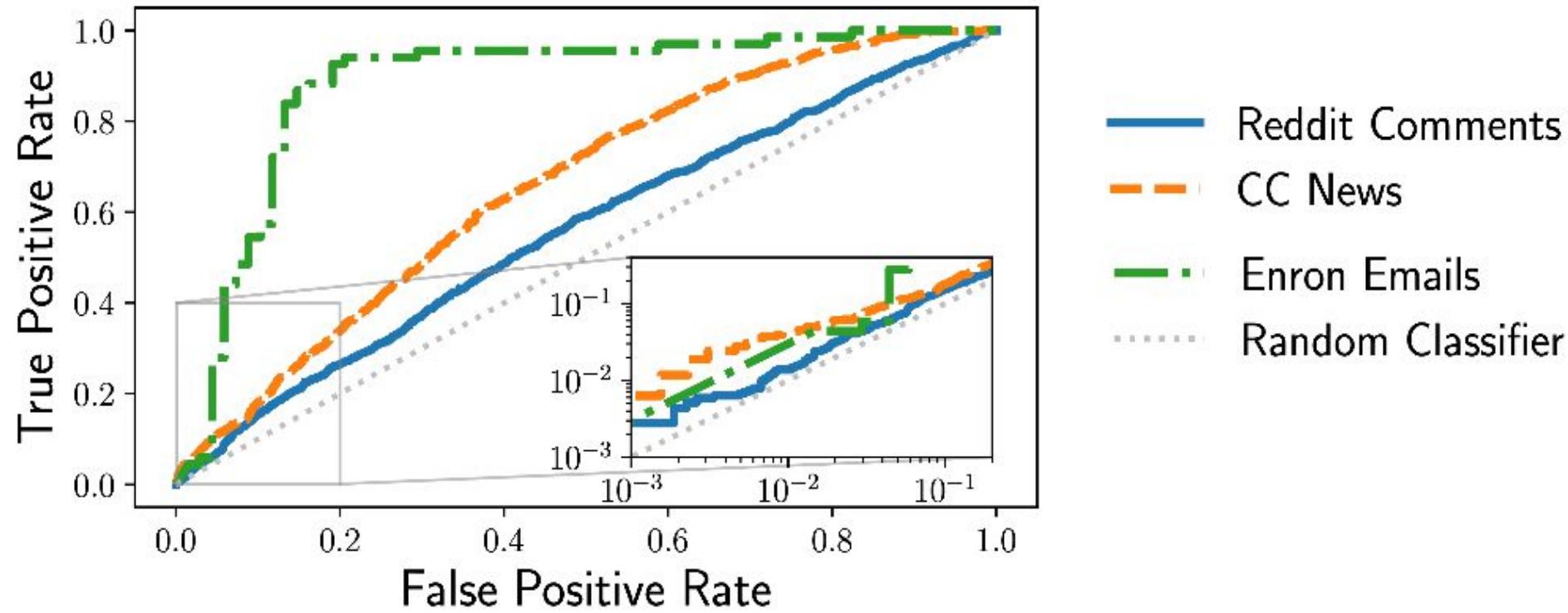
3. Test statistic

$$\hat{T}(x^{(1)}, \dots, x^{(m)}) = \frac{1}{m} \sum_{i=1}^m \log \left(\frac{p_\theta(x^{(i)})}{p_{\text{ref}}(x^{(i)})} \right)$$

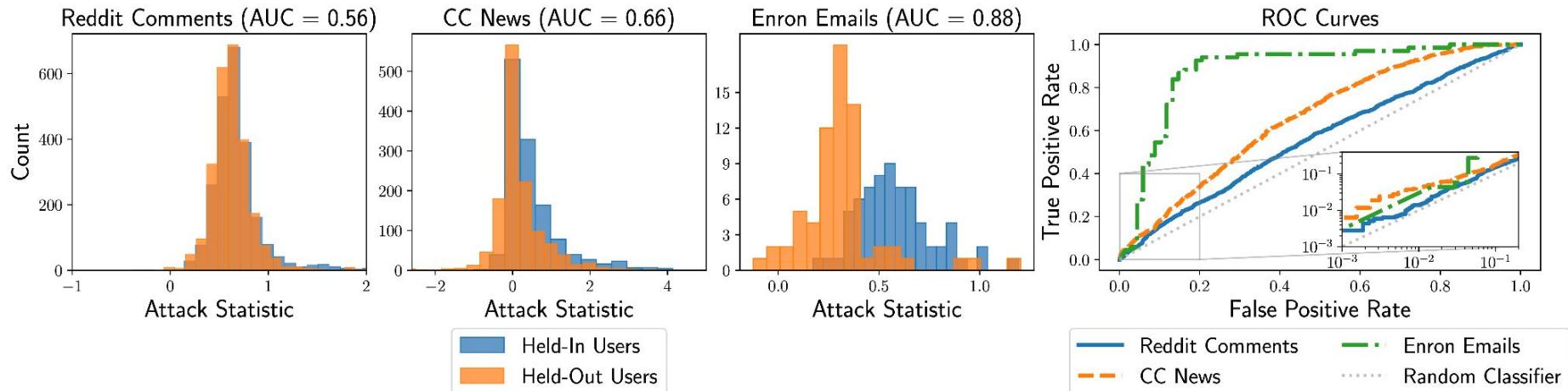
4. U was in training if $\hat{T}(x^{(1)}, \dots, x^{(m)}) > \tau$

Evaluation

ROC Curves



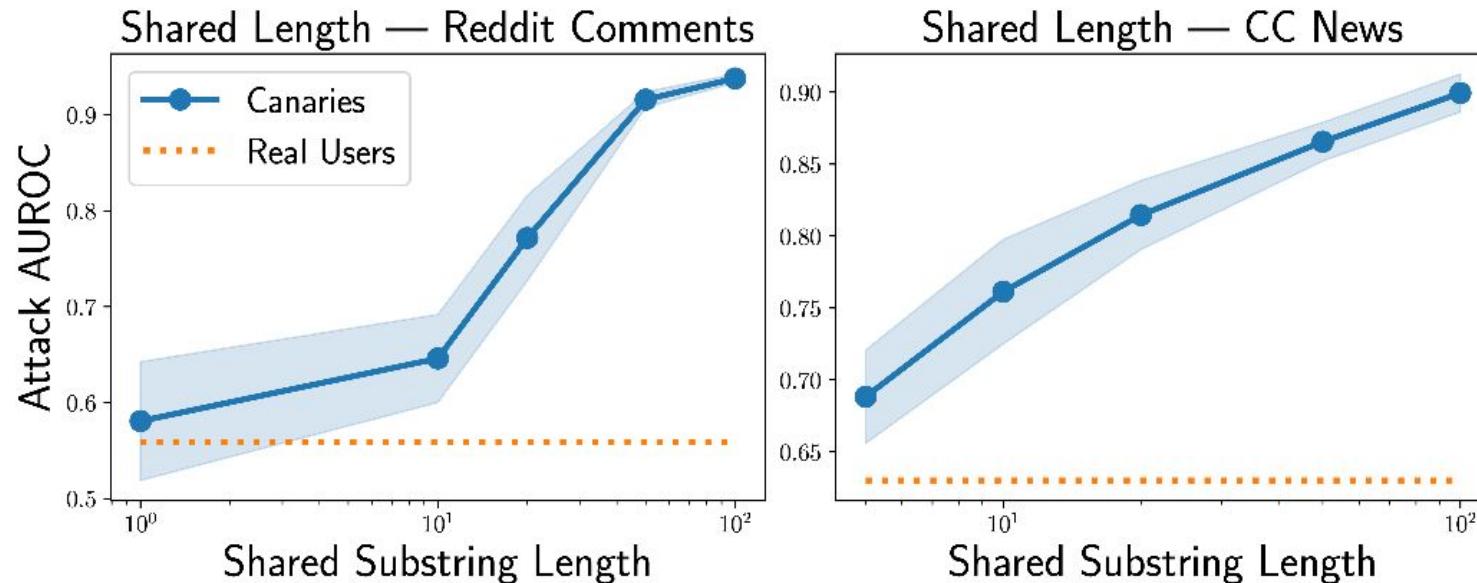
User inference is effective when #users is small and data per user is large



More fine-tuning samples per user

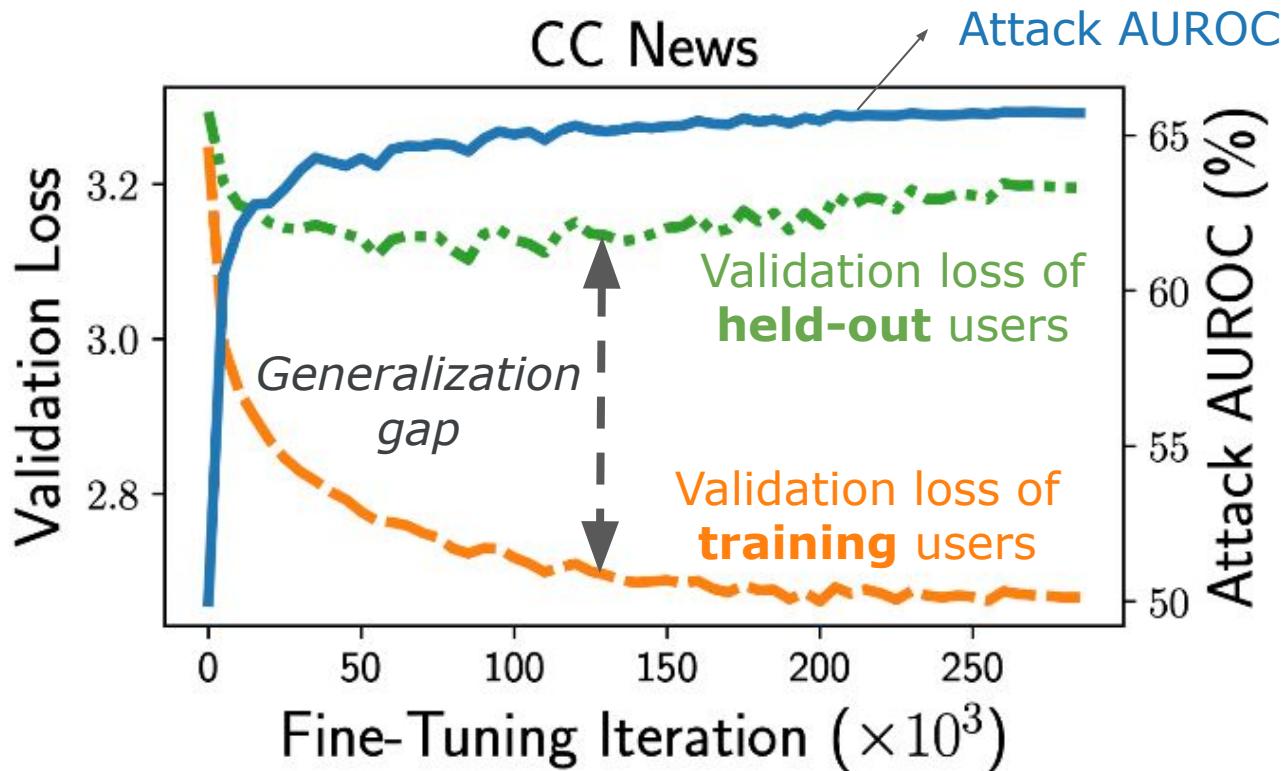
More users

Short common phrases can exacerbate user inference



User inference mechanism:

Overfitting to ***data distributions*** of training users



Spearman Correlation(Generalization gap, AUROC) = 0.995

Can user inference be mitigated?

Do not work

Early stopping

Gradient clipping

Limited Mitigation*

Data limits per user

Data deduplication

Differential privacy

Differential privacy (DP)

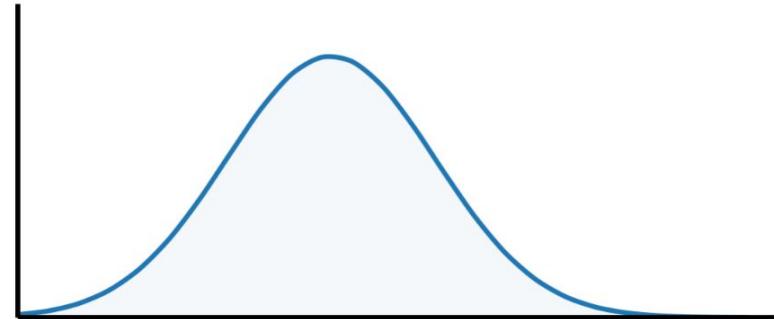
Dataset



Randomized
Algorithm

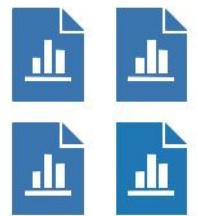


Output Distribution
(e.g. over models)



Differential privacy (DP)

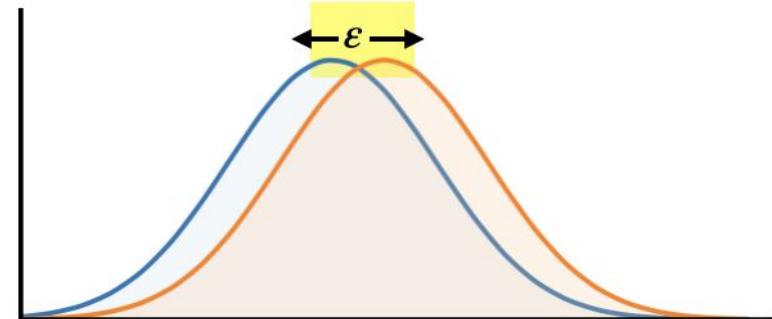
Dataset



Randomized
Algorithm



Output Distribution
(e.g. over models)



+



A randomized algorithm is **ϵ -differentially private** if the addition of **one unit of data** does not alter its output distribution by more than ϵ .

 Unit of data
= **example**

Example-level Differential privacy (DP)

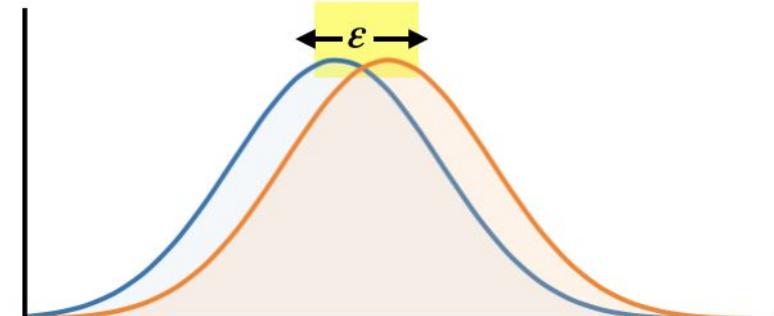
Dataset



Randomized
Algorithm



Output Distribution
(e.g. over models)



+



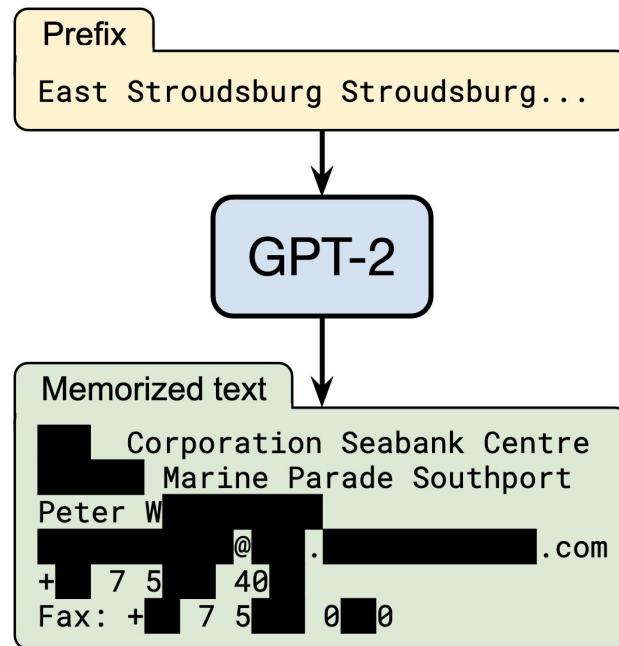
A randomized algorithm is **ε -differentially private** if the addition of **one example** does not alter its output distribution by more than ε

LONG LIVE THE REVOLUTION.
OUR NEXT MEETING WILL BE
AT| THE DOCKS AT MIDNIGHT
ON JUNE 28 TAB



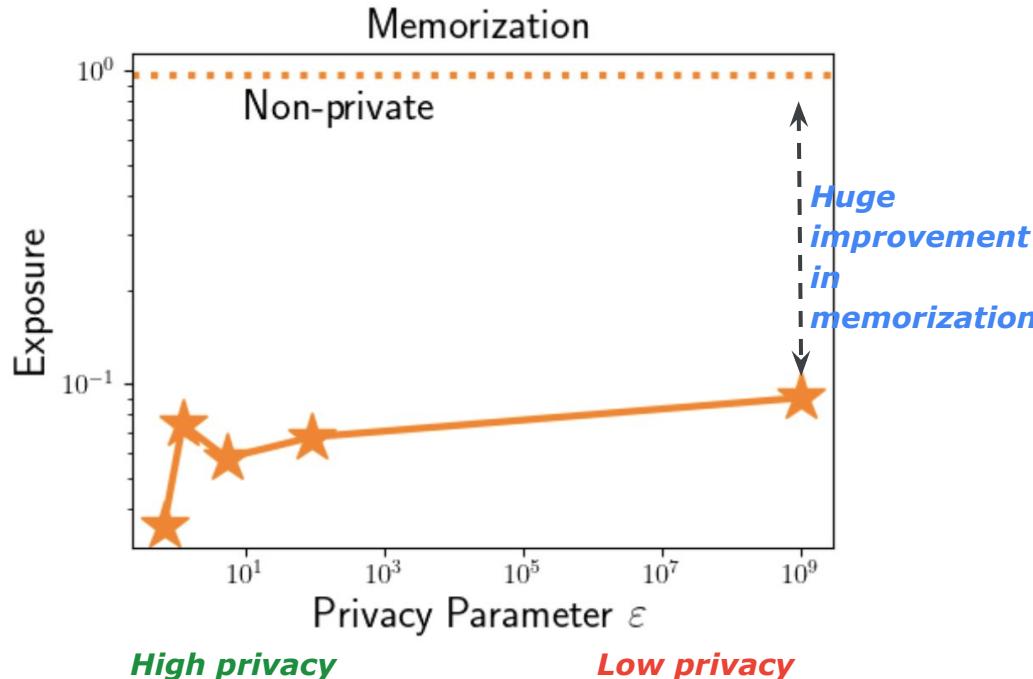
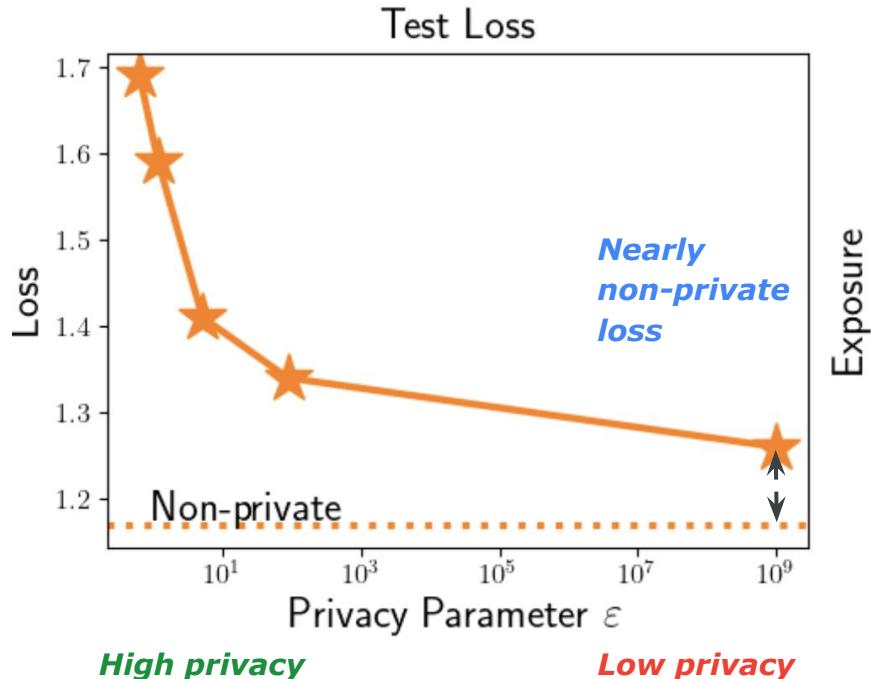
WHEN YOU TRAIN PREDICTIVE MODELS
ON INPUT FROM YOUR USERS, IT CAN
LEAK INFORMATION IN UNEXPECTED WAYS.

Models leak information about their training data



Carlini et al. (USENIX Security 2021)

Example-level DP eliminates memorization



Example-level DP offers limited mitigation for user inference

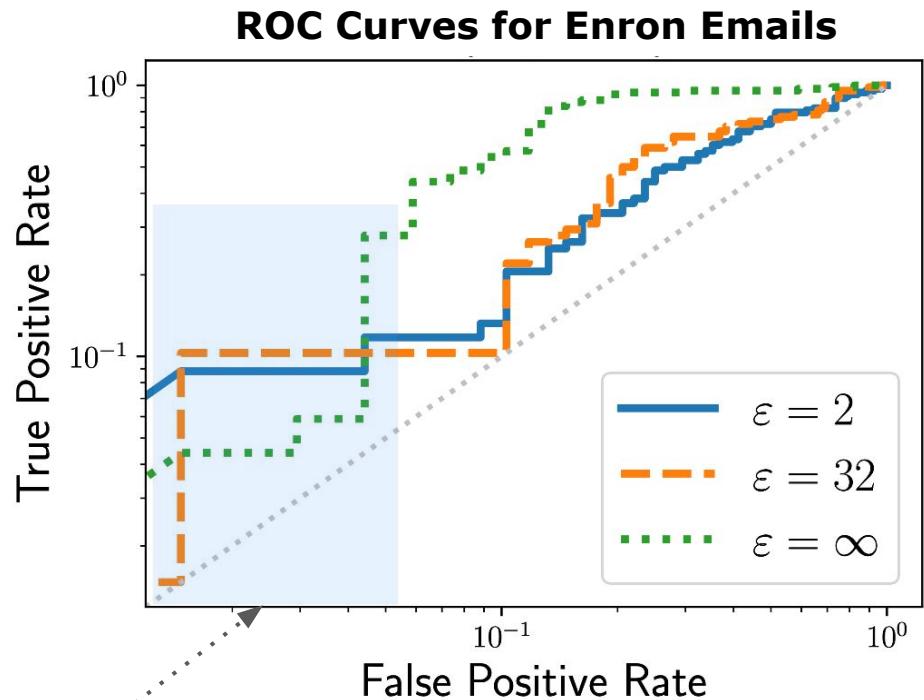
AUROC:

- non-private: 88%
- $\epsilon = 32$: 70%

Utility:

- DP model reaches what the private model achieves in 1/3 epoch

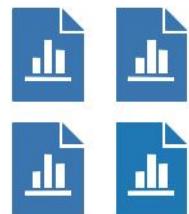
Example-level DP does not help here



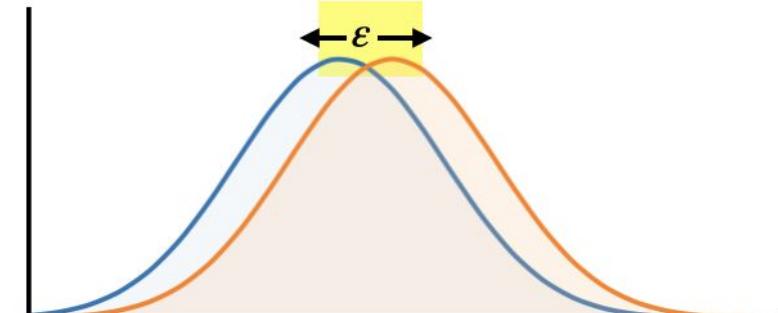
User Example-level Differential privacy (DP)

Unit of data
= user

Dataset



Output Distribution
(e.g. over models)



A randomized algorithm is **ϵ -differentially private** if the addition of **one user's data** does not alter its output distribution by more than ϵ

User-level DP: Provable protections against user inference

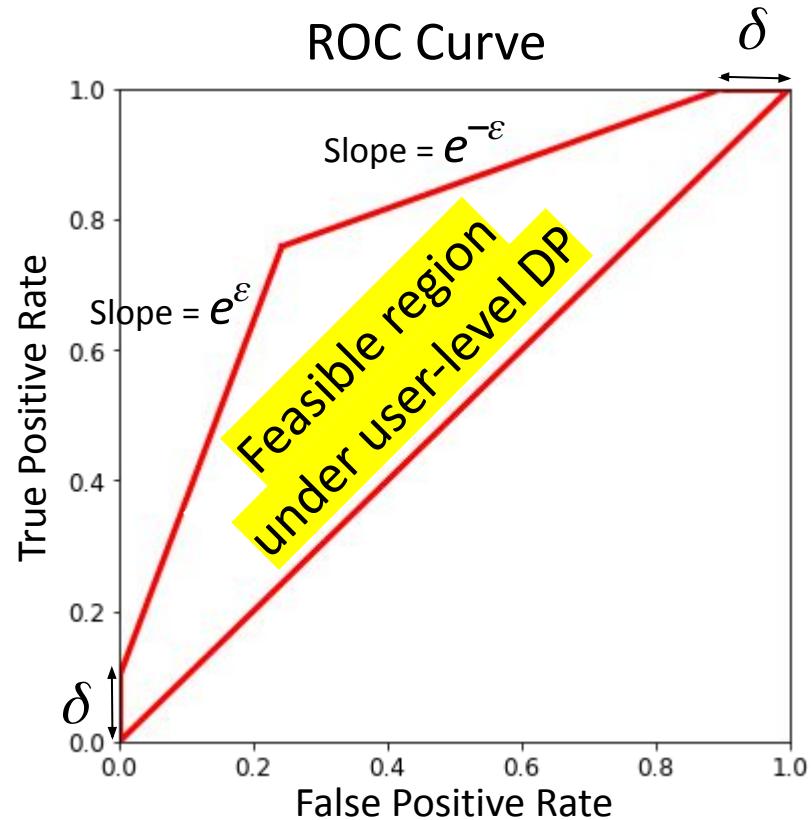
- By **differential privacy definition**:

$$\text{TPR} \leq e^\varepsilon \text{FPR} + \delta$$

True Positive Rate

False Positive Rate

- Fundamental limits on the success of membership inference



Fine-Tuning Large Language Models with User-Level Differential Privacy

Zachary Charles

Google Research

Seattle, WA, USA

zachcharles@google.com

Arun Ganesh

Google Research

Seattle, WA, USA

arunganesh@google.com

Ryan McKenna

Google Research

Seattle, WA, USA

mckennar@google.com

H. Brendan McMahan

Google Research

Seattle, WA, USA

mcmahan@google.com

Nicole Mitchell

Google Research

San Francisco, CA, USA

nicolemitchell@google.com

Krishna Pillutla

IIT Madras

Chennai, India

krishnap@dsai.iitm.ac.in

Keith Rush

Google Research

Seattle, WA, USA

krush@google.com

CORRELATED NOISE PROVABLY BEATS INDEPENDENT NOISE FOR DIFFERENTIALLY PRIVATE LEARNING

Christopher A. Choquette-Choo* Krishnamurthy (Dj) Dvijotham* Krishna Pillutla*
 Arun Ganesh Thomas Steinke Abhradeep Guha Thakurta

Efficient and Near-Optimal Noise Generation for Streaming Differential Privacy*

Krishnamurthy (Dj) Dvijotham
Google DeepMind
 dvijothamcs@gmail.com

H. Brendan McMahan
Google Research
 mcmahan@google.com

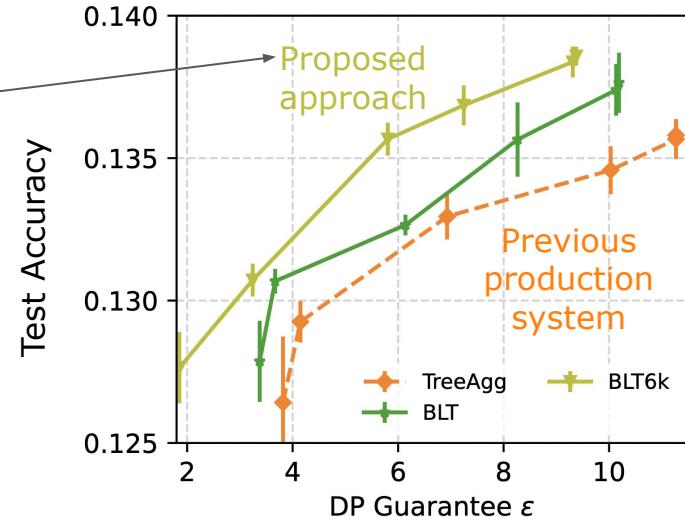
Krishna Pillutla
IIT Madras
 krishnap@dsai.iitm.ac.in

Thomas Steinke
Google DeepMind
 steinke@google.com

Abhradeep Thakurta
Google DeepMind
 athakurta@google.com

Advances in DP training

Google's production LM (Portuguese language) for next-word prediction



Plot: McMahan, Xu, Zhang (2024).

Thank you!

User Inference Attacks on Large Language Models.
EMNLP 2024 (Oral Presentation)



Nikhil Kandpal
U. Toronto



Alina Oprea
Northeastern U.



Peter Kairouz
Google



Chris Choquette
Google



Zheng Xu
Google