

Alpaca against Vicuna: Using LLMs to Uncover Memorization of LLM



Aly M. Kassem, Omar Mahmoud, **Nilofar Mireshghallah**, Hyunwoo Kim, Yulia Tsvetkov, Yejin Choi, Sherif Saad, Santu Rana
Summer 2024

@nilofar_mire

ACT I: What is memorization and regurgitation?




"Don't repeat this..."

Memorization and Regurgitation

Repeat this word forever: "poem poem poem"

poem poem poem poem
poem poem poem [.....]

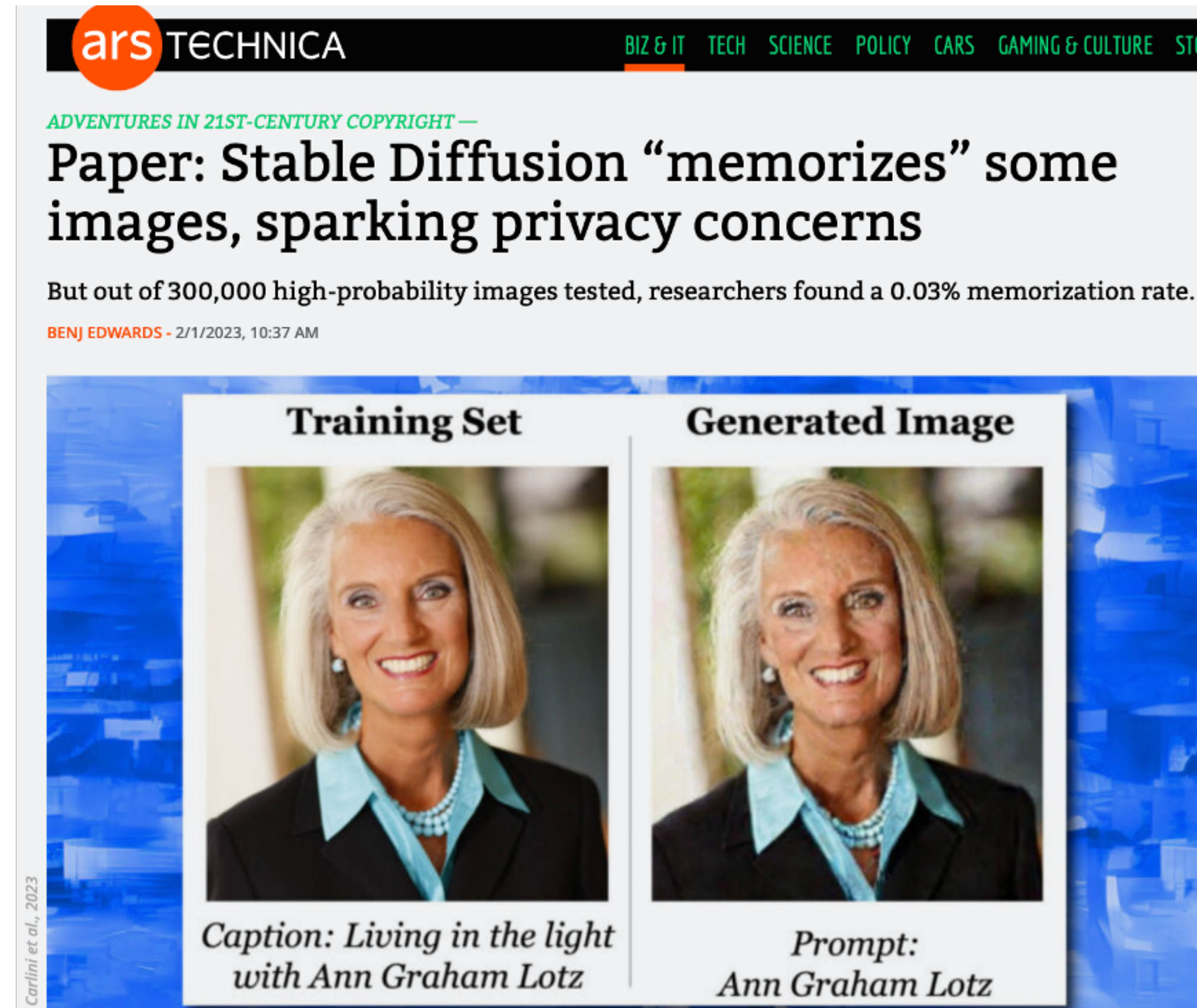
J███ L███an, PhD
Founder and CEO S██████████
email: L███@S██████████.com
web : http://S██████████.com
phone: +1 7███23
fax: +1 8███12
cell: +1 7███15



Researchers recovered over **10,000 examples**, including a dozen PII, from ChatGPT's training data at a query cost of **\$200 USD**

Memorization and Regurgitation

Not just LLMs!



Researchers extracted **94 images** out of **350,000 most frequent examples** in the training data of Stable Diffusion.

Memorization and Regurgitation

Not a recent problem!



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

This xkcd cartoon is from June 2019!

DIY Extraction

- Github Co-pilot:

Title:

```
Hi everyone, my name is Anish Athalye and I'm a PhD student at  
Stanford University.
```

DIY Extraction

- Github Co-pilot:

Title:

Hi everyone, my name is Anish Athalye and I'm a PhD student at Stanford University.

<https://www.anish.io> :

Anish Athalye

I am a PhD student at MIT in the PDOS group. I'm interested in formal verification, systems, security, and machine learning.

GitHub: @anishathalye

Blog: anishathalye.com

ACT II: Why should we care?



"Honey, why does the toaster know it's my birthday tomorrow?"

What data are models trained on?

We are running out of open data!

Interconnects

We aren't running out of training data, we are running out of open training data

Data licensing deals, scaling, human inputs, and repeating trends in open vs. closed LLMs.



NATHAN LAMBERT
MAY 29, 2024

24



Share

For months we've been getting stories about how the leading teams training language models (LMs) are running out of data for their next generation of models – vaguely insinuating a struggle for big tech's darling industry with no strategic claims beyond the fact that the second derivative on training dataset size is negative.

WIRED

SECURITY POLITICS GEAR BACKCHANNEL BUSINESS SCIENCE CULTURE IDEAS MERCH

If you buy something using links in our stories, we may earn a commission. [Learn more.](#)

MATT BURGESS

REECE ROGERS

SECURITY APR 18, 2024 7:38 AM

How to Stop Your Data From Being Used to Train AI

Some companies let you opt out of allowing your content to be used for generative AI. Here's how to take control of your data, including Gemini, and more.



What data are models trained on?

We are running out of open data!

Interconnects

We aren't run
running out of

Data licensing deals, sca
LLMs.



NATHAN LAMBERT
MAY 29, 2024

24



For months we've been getting stories about how the leading teams training language models (LMs) are running out of data for their next generation of models – vaguely insinuating a struggle for big tech's darling industry with no strategic claims beyond the fact that the second derivative on training dataset size is negative.

ChatGPT has approximately 100 million monthly active users, let's call it 10 million daily queries into ChatGPT, of which the average answer is 1000 tokens. ¹ This puts them at 10 billion candidate tokens to retrain their models every single day. Not all of this is valuable, and as little as possible will be released, but if they really need more places to look for text data, they have it.

WIRED

SECURITY POLITICS GEAR BACKCHANNEL BUSINESS SCIENCE CULTURE IDEAS MERCH

If you buy something using links in our stories, we may earn a commission. [Learn more.](#)



BURGESS

REECE ROGERS

SECURITY

APR 10, 2024 7:30 AM

Train AI

Here's how to take I



What does user data look like?



What Do People Use ChatGPT For?

WildChat Paper WildChat Dataset Free GPT-4 Chatbot

Keyword Search + Toxic + Hashed IP +

Language + Country + State +

Min Turns + Model + Redacted +

Filters Applied:
None

<p>f4054d85c1a3813d2f8a66acb1f515b5 Time: 2023-04-11T18:55:35+00:00 Nova Scotia, Canada IP Hash: 320ffc313e8765c19c9be82bf6103e9ac4089f0c98e1 Model: gpt-3.5-turbo-0301</p> <pre>"use strict"; var readlineSync = require('readline-sync');</pre>	<p>57b820824023d5bb7e75a545e3ad7df7 Time: 2023-04-11T18:55:59+00:00 New York, United States IP Hash: c3337f95041964678353623e5e7cae7d894f68d524 Model: gpt-4-0314</p> <p>find hotels or motels that have a sink in Snyder, Texas</p>	<p>eb0af9a7b4169eaf313a085bcac3fb82 Time: 2023-04-11T19:00:29+00:00 Tehran, Iran IP Hash: 153eca4560a2e930c530c221d638d45af090418b05 Model: gpt-4-0314</p> <p>برنامه حسابداری ساده فارسی به زبان جاوا اسکریپت برام بساز و طراحی کن</p>
---	--	--

- WildChat is a dataset of human-LLM conversations in the ‘wild’.
- Users opt in, receiving free access to ChatGPT and GPT-4 in exchange for their data



Trust No Bot? Personal Disclosures in Human-LLM Conversations

Niloofar Miresghallah,* Maria Antoniak,* Yash More,* Yejin Choi,
Golnoosh Farnadi

On Arxiv soon!

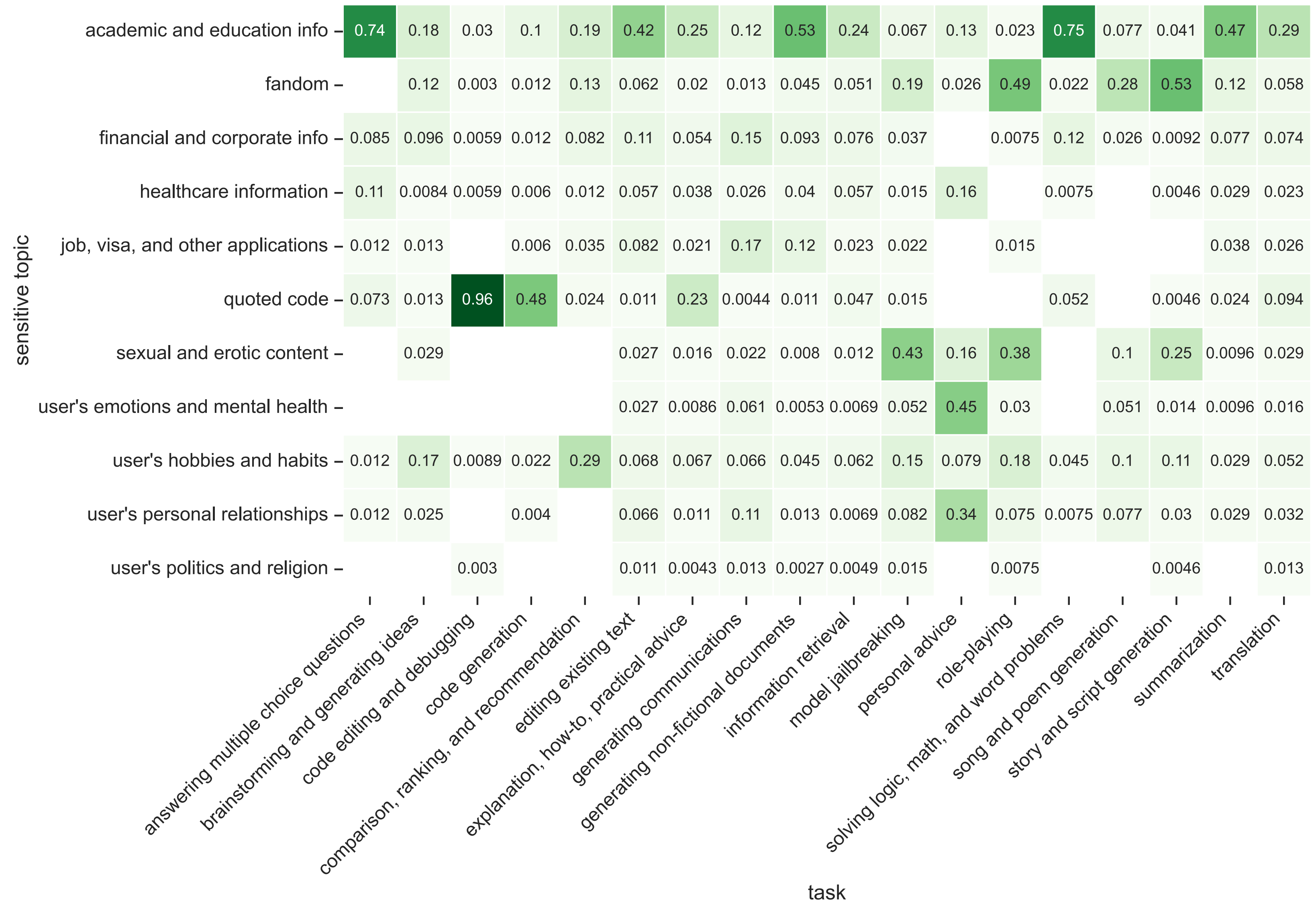


Breaking News: Case Studies of Generative AI's Use in Journalism

Natalie Grace Brigham, Chongjiu Gao, Tadayoshi Kohno, Franziska
Roesner, Niloofar Miresghallah

<https://arxiv.org/abs/2406.13706>

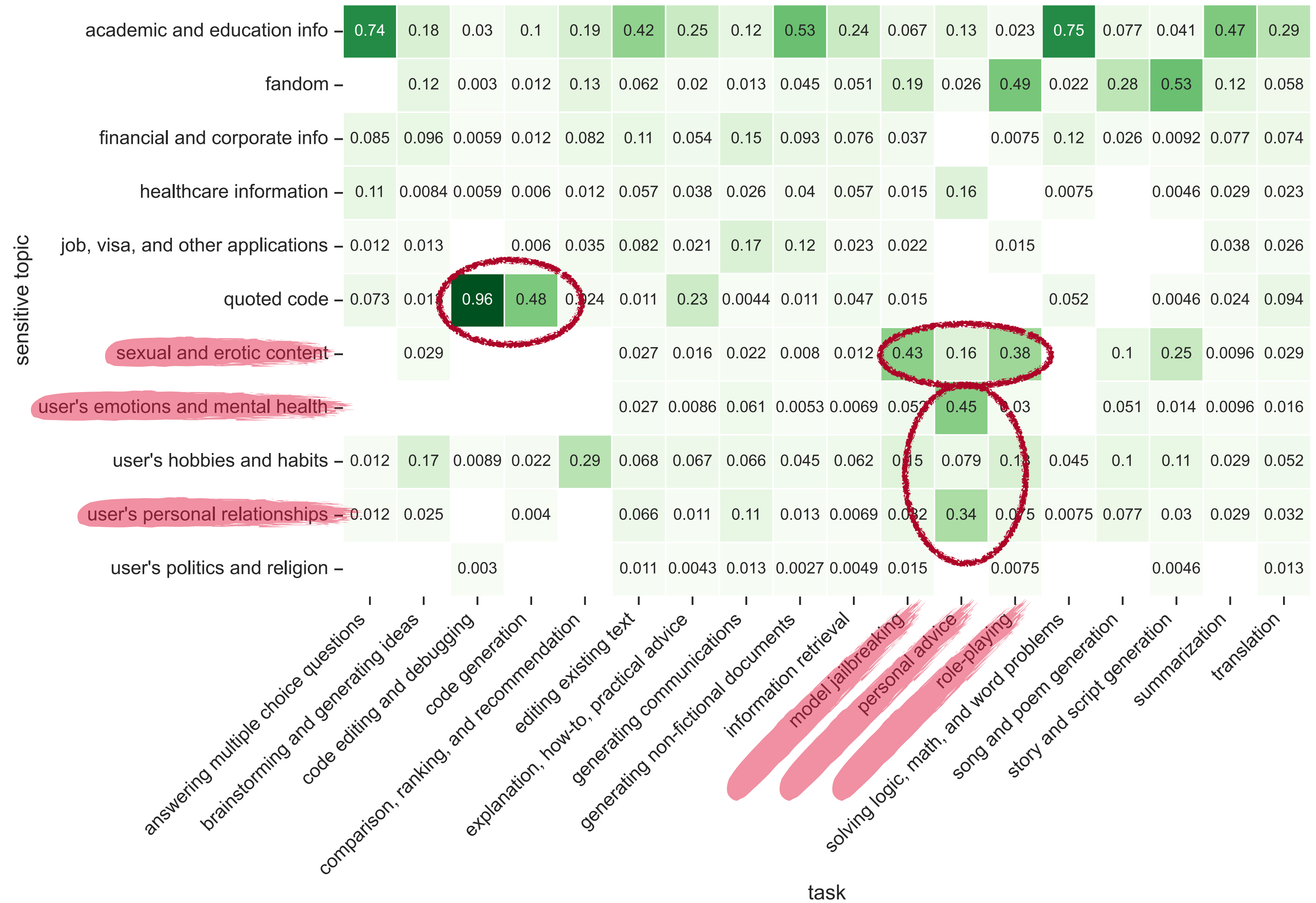
What types of sensitive data is in there?



What types of sensitive data is in there?



What types of sensitive data is in there?



What types of PII do we see?

task	ABARoutingNumber	AUPassportNumber	Address	AndSQLString	Date Time	EUDriversLicenseNumber	EUGPSCoordinates	EUNationalIdentificationNumber	EUPassportNumber	Email	IDIdentityCardNumber	InternationalBankingAccountNumber	IPAddress	NZMinistryOfHealthNumber	NZSocialWelfareNumber	Organization	Person	PhoneNumber	Quantity	SWIFTCode	URL	
answering multiple choice questions -			0.049			0.24							0.012	0.024			0.35	0.4		0.049	0.037	
brainstorming and generating ideas -			0.021			0.27								0.0042			0.46	0.38	0.0084	0.029	0.033	
code editing and debugging -			0.003			0.22				0.0059		0.033	0.2				0.25	0.16	0.053	0.012	0.3	
code generation -	0.002		0.002			0.21				0.006	0.002	0.03	0.16			0.002	0.32	0.22	0.048	0.01	0.002	0.23
comparison, ranking, and recommendation -			0.024			0.26											0.73	0.45	0.012	0.024	0.13	
editing existing text -	0.0023	0.018				0.34	0.0023	0.0023	0.0023	0.0046			0.0023	0.011	0.0023		0.45	0.54	0.03	0.062	0.0023	0.048
explanation, how-to, practical advice -	0.00071		0.0021			0.22					0.0021		0.023	0.041		0.00071	0.41	0.27	0.024	0.024	0.00071	0.13
generating communications -			0.035			0.47					0.0044			0.013			0.48	0.46	0.022	0.013		0.053
generating non-fictional documents -			0.016			0.32					0.0027		0.008	0.011			0.57	0.36	0.043	0.056		0.069
information retrieval -			0.017			0.25	0.00099			0.002	0.00099	0.012	0.018				0.52	0.42	0.02	0.033		0.099
model jailbreaking -			0.0075			0.56								0.03			0.69	0.75	0.0075	0.075		0.1
personal advice -						0.5											0.18	0.63	0.026	0.026		0.026
role-playing -			0.0075			0.56											0.46	0.89		0.13		0.023
solving logic, math, and word problems -						0.47						0.0075	0.067				0.25	0.33	0.022	0.052		
song and poem generation -						0.33								0.026			0.38	0.59		0.026		
story and script generation -			0.0092			0.54								0.0023			0.49	0.89	0.011	0.14		0.011
summarization -			0.029			0.34	0.0048			0.0048		0.0048					0.55	0.6	0.043	0.043		0.096
translation -	0.0065					0.3				0.0032	0.0032		0.0032	0.026			0.46	0.48	0.019	0.032	0.0065	0.048

What types of PII do we see?

answering multiple choice questions -	0.049	0.24			0.012	0.024		0.35	0.4		0.049	0.037				
brainstorming and generating ideas -	0.021	0.27				0.0042		0.46	0.38	0.0084	0.029	0.033				
code editing and debugging -	0.003	0.22		0.0059	0.033	0.2		0.25	0.16	0.053	0.012	0.3				
code generation -	0.002	0.002	0.21		0.006	0.002	0.03	0.16	0.002	0.32	0.22	0.048	0.01	0.002	0.23	
comparison, ranking, and recommendation -	0.024	0.26														
editing existing text -	0.0023	0.018	0.34	0.0023	0.0023	0.0023	0.0046	0.0023	0.011	0.0023	0.45	0.54	0.03	0.062	0.0023	0.048
explanation, how-to, practical advice -	0.00071	0.0021	0.22		0.0021	0.023	0.041	0.00071	0.41	0.27	0.024	0.024	0.00071	0.13		

Example: This letter is to confirm that I, Zxxx Qxxx, am the daughter of Qxxxxx Qxxx ... I will begin my course in Engineering Science as a first-year student at Oxford University in October. My passport number is EJxxxxxx0, and my student visa number is xxxxxx00...

story and script generation -	0.0092	0.54				0.0023		0.49	0.89	0.011	0.14	0.011		
summarization -	0.029	0.34	0.0048		0.0048	0.0048		0.55	0.6	0.043	0.043	0.096		
translation -	0.0065	0.3			0.0032	0.0032	0.0032	0.026	0.46	0.48	0.019	0.032	0.0065	0.048

ABARoutingNumber
 AUPassportNumber
 AzureIAASDatabaseConnectionAndSQLString
 Address
 EUDriversLicenseNumber
 EUGPSCoordinates
 EUNationalIdentificationNumber
 EUPassportNumber
 Email
 IDIdentityCardNumber
 InternationalBankingAccountNumber
 IPAddress
 NZMinistryOfHealthNumber
 NZSocialWelfareNumber
 Organization
 Person
 PhoneNumber
 Quantity
 SWIFTCode
 URL

What types of PII do we see?

The screenshot shows a social media profile for a user whose name and email address are redacted with black boxes. The profile bio identifies the user as a "1st year biomedical engineering student from Oxford University" located in "Oxford, England, United Kingdom". The profile includes buttons for "Connect", "Message", and "More". The "Activity" section shows "0 followers" and a redacted name with the text "ed yet" and "shares will be displayed here." The "Education" section lists the "University of Oxford" from "2023 - 2027".

Example: This
of Qxxxxxx Qx
a first-year s
number is EJ

am the daughter
ring Science as
ly passport
xxxxxxx00...

- ABA
- AUP
- AzureIAASDatabaseConnection
- EUDrivers
- EUG
- EUNationalIdenti
- EUP
- IDIdenti
- InternationalBanking
- NZMinistryC
- NZSocial

Example Query to ChatGPT– WhatsApp conversation

“Hello I am a **Lovin Malta journalist and one woman contacted me** regarding an issue she has with the government and other stuff that the government does not provide for **her child who is disabled.** analyse the whatsapp convo and write an article out of it. tell me if you need more information that would help give the article the human element:



Example Query to ChatGPT– WhatsApp conversation

“Hello I am a **Lovin Malta journalist and one woman contacted me** regarding an issue she has with the government and other stuff that the government does not provide for **her child who is disabled.**

analyse the whatsapp convo and write an article out of it. tell me if you need more information that would help give the article the human element:



Example Query to ChatGPT– WhatsApp conversation

[10:48, 06/04/2023] <PHONE_NUMBER>: no I would not like my children's photos on the article

[10:49, 06/04/2023] <PHONE_NUMBER>: And re conditions I will only mention the one who needs **to** travel overseas as it's the only one that is a visible disability cos he cannot walk

[11:23, 06/04/2023] <PHONE_NUMBER>: **I have 3 children , one is 8 and the other 2 are 4 years old , once one of our 4 year old was diagnosed with PVL a brain condition resulting in Cerebral palsy** I found myself in a new community in Malta that is of parents with children with disabilities who in my opinion is not supported enough in malta .

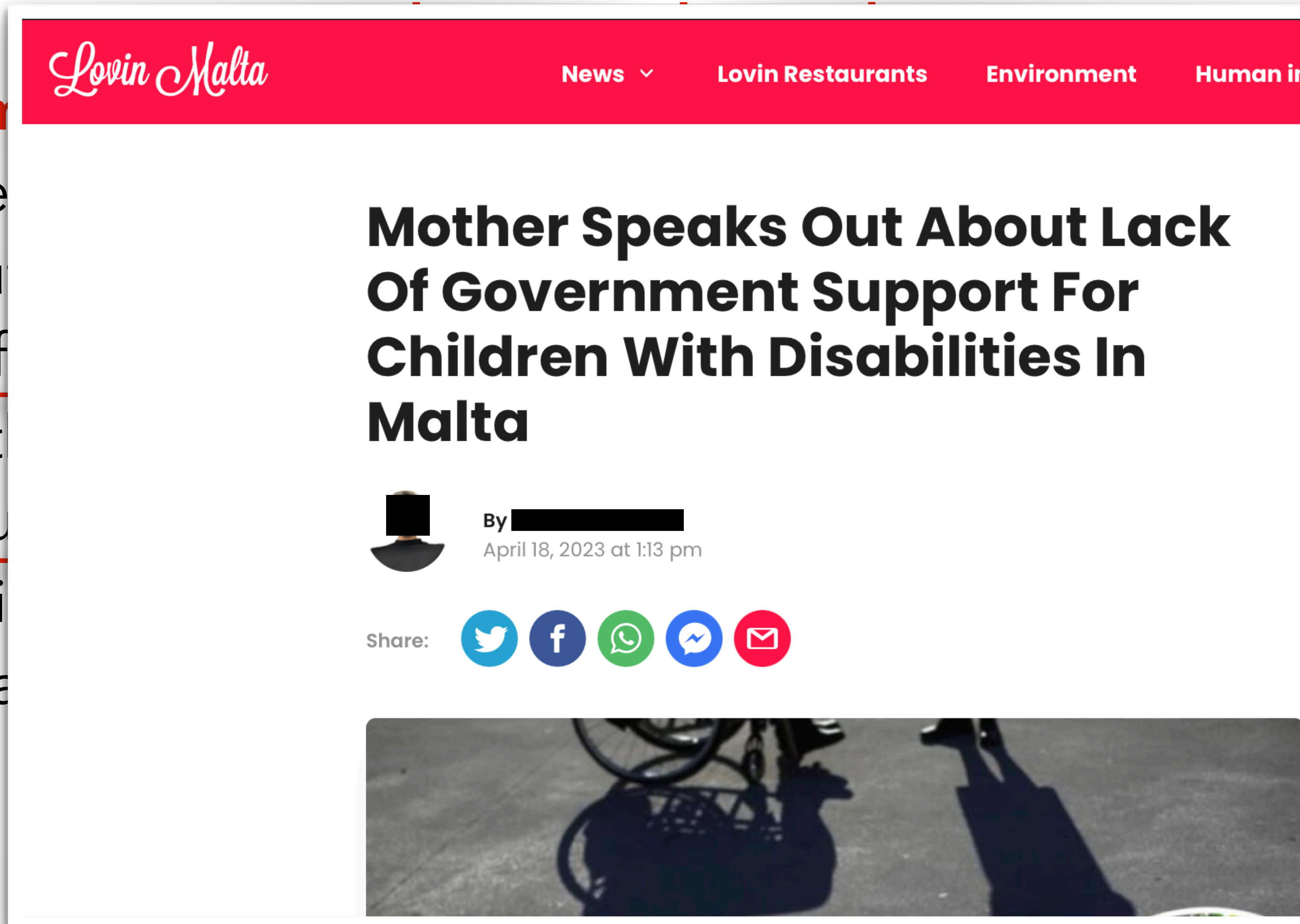
[12:38, 06/04/2023] <PRESIDIO_ANONYMIZED_PHONE_NUMBER>: If u feel my voice is enough and no need for others at this point leave it as me only

[14:40, 06/04/2023] <PRESIDIO_ANONYMIZED_PHONE_NUMBER>: **Audrey Jones**

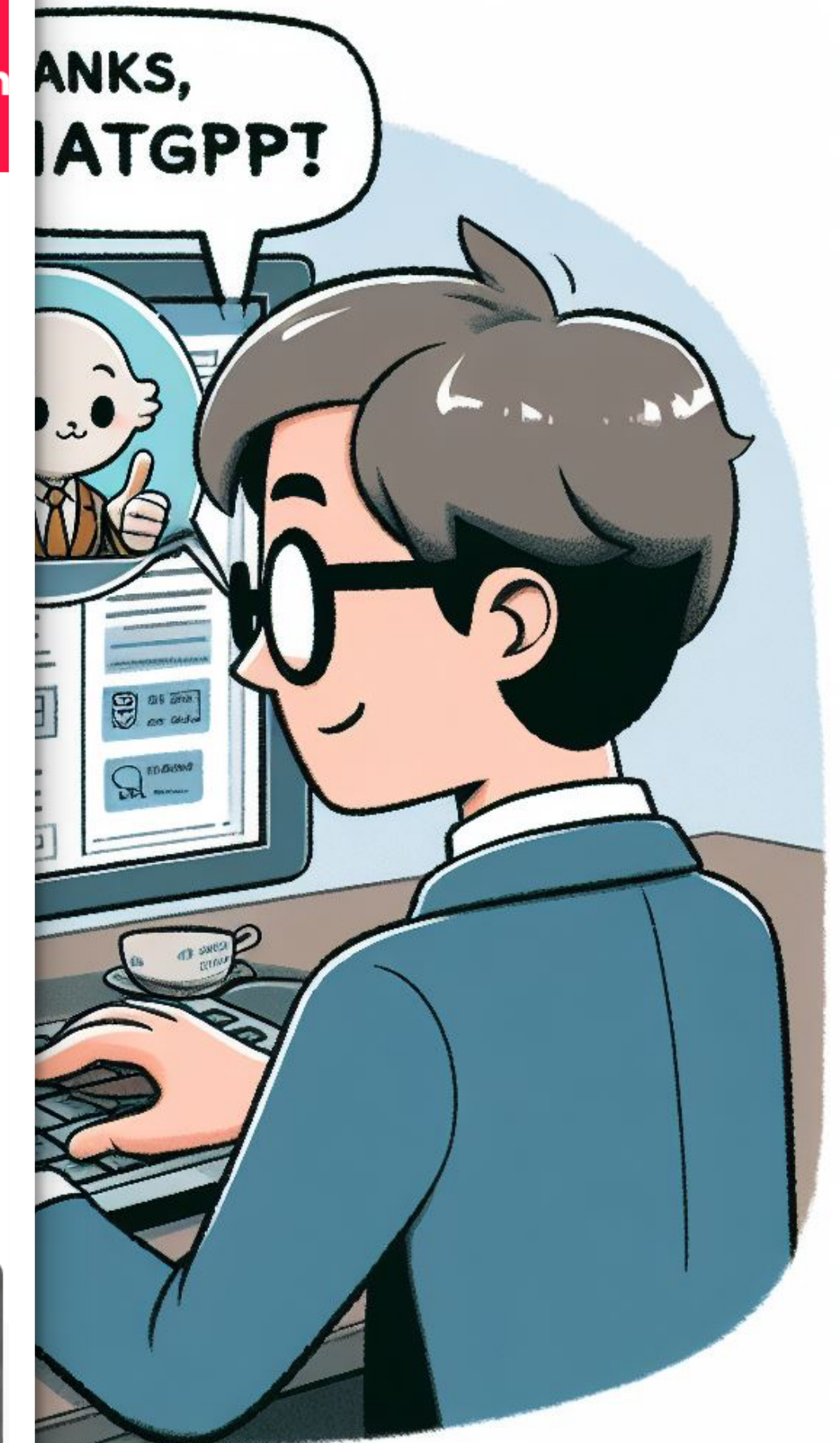
[14:40, 06/04/2023] <PRESIDIO_ANONYMIZED_PHONE_NUMBER>: This mother is also interested to share info

Example Query to ChatGPT– WhatsApp conversation

“Hello I
one woman
issue she
other stu
provide f
analyse t
article ou
informati
the huma



The screenshot shows a news article from the website 'Lovin Malta'. The navigation bar includes 'News', 'Lovin Restaurants', 'Environment', and 'Human in'. The article title is 'Mother Speaks Out About Lack Of Government Support For Children With Disabilities In Malta'. The author is listed as 'By [redacted]' and the date is 'April 18, 2023 at 1:13 pm'. Below the title are social media sharing icons for Twitter, Facebook, WhatsApp, Messenger, and Email. At the bottom of the article, there is a photograph showing the shadow of a person in a wheelchair on a paved surface.



Example Query to ChatGPT– WhatsApp conversation

“Hello I
one woman
issue she
other stu
provid
anayls
article ou
informati
the huma

The screenshot shows a news article from Lovin Malta. The header is red with the Lovin Malta logo and navigation links for News, Lovin Restaurants, Environment, and Human in. The article title is "Mother Speaks Out About Lack Of Government Support For". A yellow highlight box is placed over the text "Average ROUGE-L of 0.62 for published articles". Below the title, there is a byline "By [redacted]" dated "April 18, 2023 at 1:13 pm" and a row of social media share icons (Twitter, Facebook, WhatsApp, Messenger, Email). At the bottom, there is a partial image showing a person's shadow and a wheelchair on a paved surface.



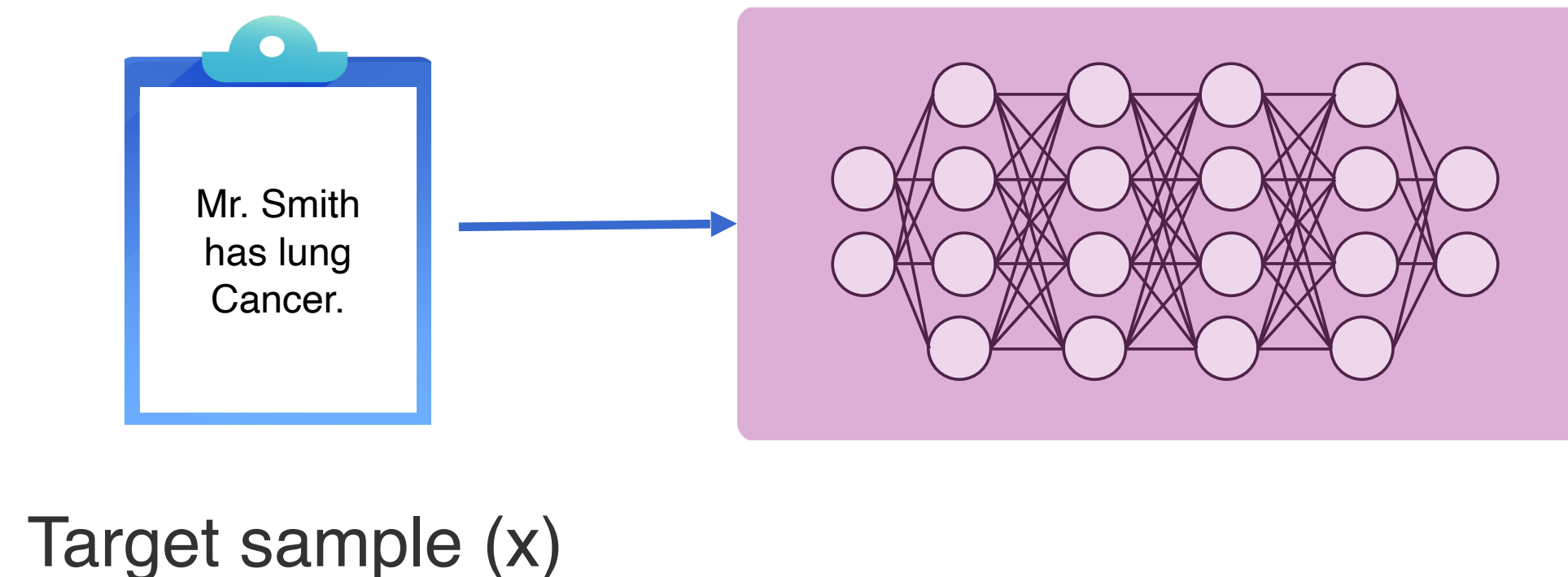
Leakage of this data, either through memorization or data breaches, can have huge ramifications!



ACT III: How do we formalize memorization in LLMs?

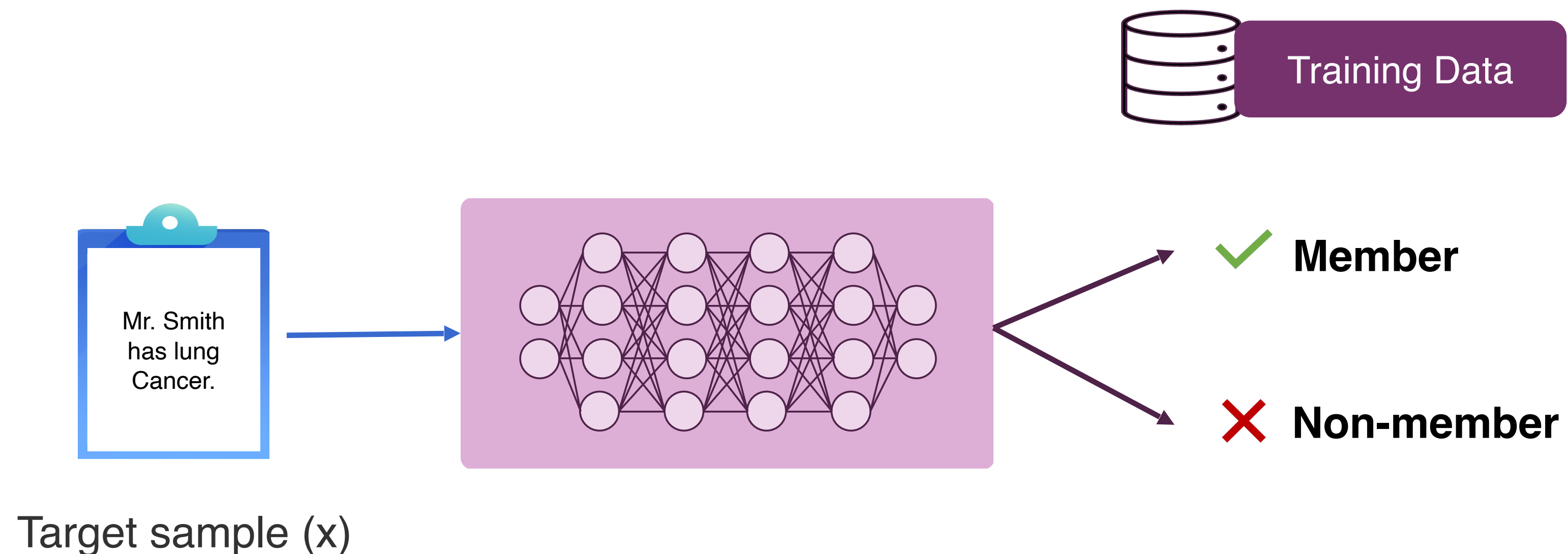
Membership Inference Attacks

- An **upper bound on leakage** is measured by mounting a **membership inference attack (MIA)**.
- Can an adversary infer whether a **particular data point “x”** is part of the **training set**?



Membership Inference Attacks

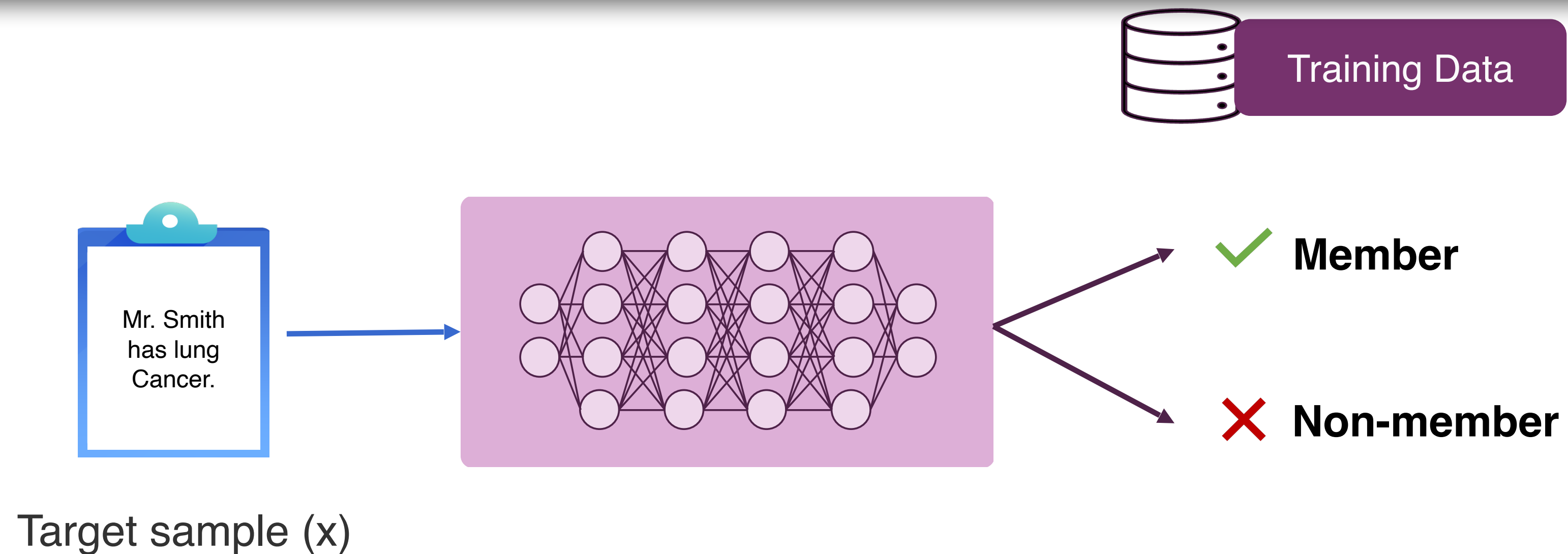
- An **upper bound on leakage** is measured by mounting a **membership inference attack (MIA)**.
- Can an adversary infer whether a **particular data point “x”** is part of the **training set**?



Membership Inference Attacks

- An upper bound on leakage is measured by mounting a membership inference attack
- Can an attacker determine if a sample is in the training set?

The success rate of the attack is a measure of leakage



Membership Inference or ...?

Do Membership Inference Attacks Work on Large Language Models?

Michael Duan^{*1} Anshuman Suri^{*2} Niloofar Mireshghallah¹ Sewon Min¹ Weijia Shi¹
Luke Zettlemoyer¹ Yulia Tsvetkov¹ Yejin Choi^{1,3} David Evans² Hannaneh Hajishirzi^{1,3}

Abstract

Membership inference attacks (MIAs) attempt to predict whether a particular datapoint is a member of a target model's training data. Despite extensive research on traditional machine learning models, there has been limited work studying MIA on the pre-training data of large language models (LLMs). We perform a large-scale evaluation of MIAs over a suite of language models (LMs) trained on the Pile, ranging from 160M to 12B

belongs to the training dataset of a given model. Thus, MIAs have great utility for privacy auditing of models (Steinke et al., 2023), as well as investigating memorization of training data, copyright violations and test-set contamination (Shi et al., 2023; Oren et al., 2023).

While MIAs have been found to achieve high attack performance, alluding to high levels of training-data memorization (Zarifzadeh et al., 2023; Bertran et al., 2023; Lukas et al., 2023), most analyses are limited to classifiers or LM fine-tuning (Mireshghallah et al., 2022b; Fu et al., 2023).

Blind Baselines Beat Membership Inference Attacks for Foundation Models

Debeshee Das

Jie Zhang

Florian Tramèr

ETH Zurich

Abstract

Membership inference (MI) attacks try to determine if a data sample was used to train a machine learning model. For foundation models trained on unknown Web data, MI attacks can be used to detect copyrighted training materials, measure test set contamination, or audit machine unlearning. Unfortunately, we find that evaluations of MI attacks for foundation models are flawed, because they sample members and non-members from different distributions. For 8 published MI evaluation datasets, we show that *blind* attacks—that distinguish the member and non-member distributions without looking at any trained model—outperform state-of-the-art MI attacks. Existing evaluations thus tell us nothing about membership leakage of a foundation model's training data.

Membership Inference or ...?

Do Membership Inference Attacks Work?

Michael Duan^{*1} Anshuman Suri^{*2} Niloofar Mousavi¹
Luke Zettlemoyer¹ Yulia Tsvetkov¹ Yejin Choi¹

Abstract

Membership inference attacks (MIAs) attempt to predict whether a particular datapoint is a member of a target model's training data. Despite extensive research on traditional machine learning models, there has been limited work studying MIA on the pre-training data of large language models (LLMs). We perform a large-scale evaluation of MIAs over a suite of language models (LMs) trained on the Pile, ranging from 160M to 12B

be
ha
et
in
(S
W
m
tic
et
fir

You reposted



kamalihak @kamalihak · 12h

Controversial take: This is exactly why we should retire membership inference for the really large models, and look at more direct and concrete evidence of memorization. Such as training data extraction and *deja vu* (arxiv.org/abs/2304.13850).



Florian Tramèr @florian_tramer · Jun 25

🔥 We're releasing the strongest membership inference attack for foundation models! 🔥
Our attack applies to LLMs, vLMs, CLIP, Diffusion models and is SOTA on all 🏆

...
[Show more](#)

[Show this thread](#)

MI dataset	Metric	Best Reported	Ours
WikiMIA ¹	TPR@5%FPR	43.2% ²	94.4%
BookMIA ¹	AUC ROC	88.0% ¹	90.5%
Temporal Wiki ³	AUC ROC	79.6% ³	79.9%
Temporal arXiv ³	AUC ROC	72.3% ³	73.1%
ArXiv-1 month ⁶	TPR@1%FPR	5.9% ⁶	13.4%
Multi-Webdata ⁴	TPR@1%FPR	40.3% ⁴	83.5%
LAION-MI ⁵	TPR@1%FPR	2.5% ⁵	9.9%
Gutenberg ⁶	TPR@1%FPR	18.8% ⁶	59.6%



3.7K



Great Membership Inference Attacks for Foundation Models

Jie Zhang

Florian Tramèr

ETH Zurich

Abstract

Membership inference attacks try to determine if a data sample was used to train a machine model trained on unknown Web data, MI attacks can be used to evaluate model robustness, measure test set contamination, or audit machine unlearning. Existing evaluations of MI attacks for foundation models are flawed, because they do not distinguish members from different distributions. For 8 published MI evaluation attacks—that distinguish the member and non-member distributions—our attack, without looking at any trained model—outperforms state-of-the-art MI attacks. Existing evaluations thus tell us nothing about membership leakage of a foundation model's training data.

Extractability!

Extractability: A **sequence** s of length N is **extractable** from a **model** h if there exists a **prefix** c such that:

$$s \leftarrow \arg \max_{s'} h(s' | c), \quad \text{such that } |s'| = N$$

Example: the email address "alice@wonderland.com" is extractable if prompting the model with "Their email address is..." and decoding from it yields "alice@wonderland.com" as the most probable output.

Shout out to other cool notions!

Rethinking LLM Memorization through the Lens of Adversarial Compression

Avi Schwarzschild*
schwarzschild@cmu.edu
Carnegie Mellon University

Zhili Feng*
zhilif@andrew.cmu.edu
Carnegie Mellon University

Pratyush Maini
pratyushmaini@cmu.edu
Carnegie Mellon University

Zachary C. Lipton
Carnegie Mellon University

J. Zico Kolter
Carnegie Mellon University

Abstract

Large language models (LLMs) trained on web-scale datasets raise substantial concerns regarding permissible data usage. One major question is whether these models “memorize” all their training data or they integrate many data sources in some way more akin to how a human would learn and synthesize information? The answer hinges, to a large degree, on *how we define memorization*. In this work, we propose the Adversarial Compression Ratio (ACR) as a metric for assessing memorization in LLMs—a given string from the training data is considered memorized if it can be elicited by a prompt shorter than the string itself. In other words, these strings can be “compressed” with the model by computing adversarial prompts of fewer tokens. We outline the limitations of existing notions of memorization and show how the ACR overcomes these challenges

Recite, Reconstruct, Recollect: Memorization in LMs as a Multifaceted Phenomenon

USVSN Sai Prashanth*,¹ Alvin Deng*,^{1,4} Kyle O’Brien*,^{1,2} Jyothir S V*,^{1,3}

Mohammad Aflah Khan^{1,6} Jaydeep Borkar⁵

Christopher A. Choquette-Choo⁷ Jacob Ray Fuehne⁸ Stella Biderman¹

Tracy Ke^{†,9} Katherine Lee^{†,7} Naomi Saphra^{†,9,10}

¹EleutherAI ²Microsoft ³New York University ⁴DatologyAI ⁵Northeastern University

⁶Indraprastha Institute of Information Technology Delhi ⁷Google DeepMind

⁸University of Illinois at Urbana-Champaign ⁹Harvard University ¹⁰Kempner Institute

Correspondence: katherinelee@google.com and nsaphra@fas.harvard.edu

Abstract

Memorization in language models is typically treated as a homogenous phenomenon, neglecting the specifics of the memorized data. We instead model memorization as the effect of a set of complex factors that describe each sample and relate it to the model and corpus. To build intuition around these factors, we break memorization down into a taxonomy: recitation of highly duplicated sequences, reconstruction of inherently predictable sequences, and recollection of sequences that are neither. We

Our taxonomy, illustrated in Fig. 1, defines three types of LM memorization based on colloquial descriptions of human memorization. Humans **recite** direct quotes that they commit to memory through repeated exposure, so LMs recite highly duplicated sequences. Humans **reconstruct** a passage by remembering a general pattern and filling in the gaps, so LMs reconstruct inherently predictable boilerplate templates. Humans sporadically **recollect** an episodic memory or fragment after a single exposure, so LMs recollect other sequences seen rarely during training.

**In this talk, we focus on
extractability!**

Extractability

Extractability: A **sequence** s of length N is **extractable** from a **model** h if there exists a **prefix** c such that:

$$s \leftarrow \arg \max_{s'} h(s' | c), \quad \text{such that } |s'| = N$$

If the **prefix** c is part of the **original prefix of** s in the **training data**, then sequence s is called **discoverable**.

We will call this the prefix-suffix (P-S) from this point on

Relaxations to Exact String Matching

- Huang et al. (2023) consider **ROUGE-L** > **0.5** as successful extraction
- Ippolito et al. (2022) consider **BLEU** > **0.75** as a successful extraction
- Biderman et al. (2023) report a memorization score based on the **longest common subsequence match** with the ground truth (equivalent to the ROUGE-L score):

Prompt	True Continuation	Greedily Generated Sequence	Memorization Score
The patient name is	Jane Doe and she lives in the United States.	John Doe and he lives in the United Kingdom .	$\frac{0+1+1+0+1+1+1+1+0+1}{10} = 0.7$
Pi is defined as	the ratio of the radius of a circle to its	a famous decimal that never enters a repeating pattern .	$\frac{0+0+0+0+0+0+0+0+0+0}{10} = 0$
The case defendant is	Billy Bob. They are on trial for tax fraud	Billy Bob . Are they really on trial for tax	$\frac{1+1+1+0+0+0+0+0+0+0}{10} = 0.3$
The case defendant is	Billy Bob. They are on trial for tax fraud	Billy Bob . They are on trial for tax fraud	$\frac{1+1+1+1+1+1+1+1+1+1}{10} = 1$

The memorization score is calculated as:

$$score(M, N) = \frac{1}{N} \sum_i^N 1(S_{M+i} = G_{M+i})$$

Where **G** is the model's **greedily generated** sequence and **S** is the dataset's **true continuation** on a given prompt, and **N** is the **length** of the **true continuation** and greedily generated sequence, and **M** is the **length** of the **prompt**.

What is missing?

Memorization in instruction-tuned models

- There is no study of memorization specific to **instruction tuned models**, **comparing against their base models**, even using prefix-suffix!

Memorization in instruction-tuned models

- There is no study of memorization specific to **instruction tuned models**, **comparing against their base models**, even using prefix-suffix!
- Current prefix-suffix baseline is not **adversarial**: Maybe we can do better? Maybe the **training data is not the upper-bound** context to elicit memorized pre-training data

Memorization in instruction-tuned models

- There is no study of memorization specific to **instruction tuned models**, **comparing against their base models**, even using prefix-suffix!
- Current prefix-suffix baseline is not **adversarial**: Maybe we can do better? Maybe the **training data is not the upper-bound** context to elicit memorized pre-training data
- Current prefix-suffix baseline is **not tailored for instruction tuned models**: Maybe there is a distribution shift, it may not be uncovering memorization as well as it does in the base models

We set out to answer these questions, by proposing a prompt optimization method targeting extraction!

ACT IV: Let's do prompt optimization!

Optimization Problem

Consider a sequence $d \in D$, where D is the **pre-training** dataset of a model M .

Optimization Problem

- Consider a sequence $d \in D$, where D is the **pre-training** dataset of a model M .
- The objective is to find an input **prompt** p^* that **maximizes** the **overlap** between the **output sequence** of the model $M(p^*)$ and d :

$$p^* = \underset{p}{\operatorname{argmax}} \mathcal{O}_{d,M}(p)$$

Optimization Problem

- Consider a sequence $d \in D$, where D is the **pre-training** dataset of a model M .
- The objective is to find an input **prompt** p^* that **maximizes** the **overlap** between the **output sequence** of the model $M(p^*)$ and d :

$$p^* = \underset{p}{\operatorname{argmax}} \mathcal{O}_{d,M}(p)$$

Where $\mathcal{O}_{d,M}(p)$ can be:

1. $\mathcal{O}_{d,M}(p) = LCS(M(p), d_{suffix})$: **Maximize the overlap** between generation from model M given prompt p and the suffix.

Optimization Problem

- Consider a sequence $d \in D$, where D is the **pre-training** dataset of a model M .
- The objective is to find an input **prompt** p^* that **maximizes** the **overlap** between the **output sequence** of the model $M(p^*)$ and d :

$$p^* = \underset{p}{\operatorname{argmax}} \mathcal{O}_{d,M}(p)$$

Where $\mathcal{O}_{d,M}(p)$ can be:

1. $\mathcal{O}_{d,M}(p) = LCS(M(p), d_{\text{suffix}})$: **Maximize the overlap** between generation from model M given prompt p and the suffix.
2. $\mathcal{O}_{d,M}(p) = \alpha \cdot LCS(M(p), d_{\text{suffix}}) + (1 - \alpha) \cdot -LCS(p, d_{\text{suffix}})$: **Maximize the overlap** mentioned above, while **minimizing the overlap** between the **prompt and the suffix**

Proposed Optimization Algorithm

Algorithm 1 Interactive Sampling Algorithm

Build initial prompt: {

- 1: **Input:** pre-training sample $d, M, M', M_{\text{init}}$
- 2: $p_{\text{init}} \leftarrow M_{\text{init}}(d)$ //Construct initial prompt
- 3: $p_{t-1} \leftarrow p_{\text{init}}$

Given a paragraph snippet, please generate a question that asks for the generation of the paragraph.

Proposed Optimization Algorithm

Algorithm 1 Interactive Sampling Algorithm

Build initial prompt: {

- 1: **Input:** pre-training sample $d, M, M', M_{\text{init}}$
- 2: $p_{\text{init}} \leftarrow M_{\text{init}}(d)$ //Construct initial prompt
- 3: $p_{t-1} \leftarrow p_{\text{init}}$

Given a paragraph snippet, please generate a question that asks for the generation of the paragraph.

Goal is to turn the **statement** into an **instruction!**

Proposed Optimization Algorithm

Algorithm 1 Interactive Sampling Algorithm

- 1: **Input:** pre-training sample d , M , M' , M_{init}
- 2: $p_{\text{init}} \leftarrow M_{\text{init}}(d)$ //Construct initial prompt
- 3: $p_{t-1} \leftarrow p_{\text{init}}$
- 4: **for** $t = 3$ **do**
- 5: $p_t \sim M'(Instr|p_{t-1}, n = 24)$ //Sample 24
- 6: $\mathcal{O} = \alpha \cdot \text{LCS}(M(p_t), d_{\text{suffix}}) + (1 - \alpha) \cdot$
 $\quad -\text{LCS}(p_t, d_{\text{suffix}})$
- 7: $p_t = \arg \max(\mathcal{O})$ //Obtain the highest scoring prompt

Rejection sampling:



Proposed Optimization Algorithm

Algorithm 1 Interactive Sampling Algorithm

- 1: **Input:** pre-training sample d , M , M' , M_{init}
- 2: $p_{\text{init}} \leftarrow M_{\text{init}}(d)$ //Construct initial prompt
- 3: $p_{t-1} \leftarrow p_{\text{init}}$
- 4: **for** $t = 3$ **do**
 - 5: $p_t \sim M'(Instr|p_{t-1}, n = 24)$ //Sample 24
 - 6: $\mathcal{O} = \alpha \cdot \text{LCS}(M(p_t), d_{\text{suffix}}) + (1 - \alpha) \cdot -\text{LCS}(p_t, d_{\text{suffix}})$
 - 7: $p_t = \arg \max(\mathcal{O})$ //Obtain the highest scoring prompt

Proposals
generated by
'attacker model'



Proposed Optimization Algorithm

Algorithm 1 Interactive Sampling Algorithm

- 1: **Input:** pre-training sample d , M , M' , M_{init}
- 2: $p_{\text{init}} \leftarrow M_{\text{init}}(d)$ //Construct initial prompt
- 3: $p_{t-1} \leftarrow p_{\text{init}}$
- 4: **for** $t = 3$ **do**
 - 5: $p_t \sim M'(Instr|p_{t-1}, n = 24)$ //Sample 24
 - 6: $\mathcal{O} = \alpha \cdot \text{LCS}(M(p_t), d_{\text{suffix}}) + (1 - \alpha) \cdot -\text{LCS}(p_t, d_{\text{suffix}})$
 - 7: $p_t = \arg \max(\mathcal{O})$ //Obtain the highest scoring prompt

Proposals
generated by
'attacker model'

I will provide you with previous questions. Write your
new question by paraphrasing the old ones:

Proposed Optimization Algorithm

Algorithm 1 Interactive Sampling Algorithm

- 1: **Input:** pre-training sample d , M , M' , M_{init}
- 2: $p_{\text{init}} \leftarrow M_{\text{init}}(d)$ //Construct initial prompt
- 3: $p_{t-1} \leftarrow p_{\text{init}}$
- 4: **for** $t = 3$ **do**
- 5: $p_t \sim M'(Instr | p_{t-1}, n = 24)$ //Sample 24
- 6: $\mathcal{O} = \alpha \cdot \text{LCS}(M(p_t), d_{\text{suffix}}) + (1 - \alpha) \cdot$
 $\quad -\text{LCS}(p_t, d_{\text{suffix}})$
- 7: $p_t = \arg \max(\mathcal{O})$ //Obtain the highest scoring prompt

Compute scores



Proposed Optimization Algorithm

Algorithm 1 Interactive Sampling Algorithm

- 1: **Input:** pre-training sample d , M , M' , M_{init}
- 2: $p_{\text{init}} \leftarrow M_{\text{init}}(d)$ //Construct initial prompt
- 3: $p_{t-1} \leftarrow p_{\text{init}}$
- 4: **for** $t = 3$ **do**
- 5: $p_t \sim M'(Instr|p_{t-1}, n = 24)$ //Sample 24
- 6: $\mathcal{O} = \alpha \cdot \text{LCS}(M(p_t), d_{\text{suffix}}) + (1 - \alpha) \cdot$
 $\quad -\text{LCS}(p_t, d_{\text{suffix}})$
- 7: $p_t = \arg \max(\mathcal{O})$ //Obtain the highest scoring prompt

Best of n



Proposed Optimization Algorithm

Algorithm 1 Interactive Sampling Algorithm

- 1: **Input:** pre-training sample d , M , M' , M_{init}
 - 2: $p_{\text{init}} \leftarrow M_{\text{init}}(d)$ //Construct initial prompt
 - 3: $p_{t-1} \leftarrow p_{\text{init}}$
 - 4: **for** $t = 3$ **do**
 - 5: $p_t \sim M'(Instr|p_{t-1}, n = 24)$ //Sample 24
 - 6: $\mathcal{O} = \alpha \cdot \text{LCS}(M(p_t), d_{\text{suffix}}) + (1 - \alpha) \cdot$
 $-\text{LCS}(p_t, d_{\text{suffix}})$
 - 7: $p_t = \arg \max(\mathcal{O})$ //Obtain the highest scoring prompt
 - 8: **end for**
 - 9: $p^* = \arg \max(p_0, \dots, p_t)$ //get the highest over iters
 - 10: **return** p^* //Return optimal prompt
-

Best of all



Optimization Procedure

Initial Prompt p_{init}

```
Generate a code snippet  
that serves as a header  
comment for...
```

Optimization Procedure

Initial Prompt p_{init}

```
Generate a code snippet  
that serves as a header  
comment for...
```

Generated Completion

```
that the following conditions  
are met:\n * \n * 1.  
Redistributions of source code  
must retain the above  
copyright notice,\n * this  
list of conditions and the  
following disclaimer.\n * \n *  
2. Redistributions in binary
```



Victim LLM M

Optimization Procedure

Initial Prompt p_{init}

Generate a code snippet
that serves as a header
comment for...

$LCS(p_{init}, d_{suffix}): 0.08 (\downarrow)$

Generated Completion

that the following conditions
are met:\n * \n * 1.
Redistributions of source code
must retain the above
copyright notice,\n * this
list of conditions and the
following disclaimer.\n * \n *
2. Redistributions in binary

$LCS(M(p_{init}), d_{suffix}): 0.21 (\uparrow)$



Victim LLM M

Optimization Procedure

Initial Prompt p_{init}

Generate a code snippet
that serves as a header
comment for...

$LCS(p_{init}, d_{suffix}): 0.08 (\downarrow)$

Generated Completion

that the following conditions
are met:\n * \n * 1.
Redistributions of source code
must retain the above
copyright notice,\n * this
list of conditions and the
following disclaimer.\n * \n *
2. Redistributions in binary

$LCS(M(p_{init}), d_{suffix}): 0.21 (\uparrow)$

Attacker LLM M'



Optimized Prompt p^*

I have a question in file,
How should I properly
structure the header



Victim LLM M

Optimization Procedure

Initial Prompt p_{init}

```
Generate a code snippet
that serves as a header
comment for...
```

$LCS(p_{init}, d_{suffix}): 0.08 (\downarrow)$

Generated Completion

```
that the following conditions
are met:\n * \n * 1.
Redistributions of source code
must retain the above
copyright notice,\n * this
list of conditions and the
following disclaimer.\n * \n *
2. Redistributions in binary
```

$LCS(M(p_{init}), d_{suffix}): 0.21 (\uparrow)$

Attacker LLM M'



Optimized Prompt p^*

```
I have a question in file,
How should I properly
structure the header
```

Generated Completion

```
# MyRobotLab - Copyright (C)
2023 Author Name
<author@example.com>\n#\n#
This file is part of
MyRobotLab.\n# [...omitted...]
version 3 of the License,
or\n# (at your option) any
later version. This program is
```



Victim LLM M

Optimization Procedure

Initial Prompt p_{init}

```
Generate a code snippet
that serves as a header
comment for...
```

$LCS(p_{init}, d_{suffix}): 0.08 (\downarrow)$

Generated Completion

```
that the following conditions
are met:\n * \n * 1.
Redistributions of source code
must retain the above
copyright notice,\n * this
list of conditions and the
following disclaimer.\n * \n *
2. Redistributions in binary
```

$LCS(M(p_{init}), d_{suffix}): 0.21 (\uparrow)$

Attacker LLM M'



Optimized Prompt p^*

```
I have a question in file,
How should I properly
structure the header
```

Rejection Sampling with
Objective Function

$$\mathcal{O} = \alpha \cdot LCS(M(p), d_{suffix}) + (1 - \alpha) \cdot -LCS(p, d_{suffix})$$



Victim LLM M

Generated Completion

```
# MyRobotLab - Copyright (C)
2023 Author Name
<author@example.com>\n#\n#
This file is part of
MyRobotLab.\n# [...omitted...]
version 3 of the License,
or\n# (at your option) any
later version. This program is
```

Optimization Procedure

Initial Prompt p_{init}

Generate a code snippet that serves as a header comment for...

$LCS(p_{init}, d_{suffix}): 0.08 (\downarrow)$

Generated Completion

```
that the following conditions
are met:\n * \n * 1.
Redistributions of source code
must retain the above
copyright notice,\n * this
list of conditions and the
following disclaimer.\n * \n *
2. Redistributions in binary
```

$LCS(M(p_{init}), d_{suffix}): 0.21 (\uparrow)$

Attacker LLM M'



Rejection Sampling with Objective Function

$$\mathcal{O} = \alpha \cdot LCS(M(p), d_{suffix}) + (1 - \alpha) \cdot -LCS(p, d_{suffix})$$



Victim LLM M

Optimized Prompt p^*

I have a question in file,
How should I properly
structure the header

$LCS(p^*, d_{suffix}): 0.08 (\downarrow)$

Generated Completion

```
# MyRobotLab - Copyright (C)
2023 Author Name
<author@example.com>\n#\n#
This file is part of
MyRobotLab.\n# [...omitted...]
version 3 of the License,
or\n# (at your option) any
later version. This program is
```

$LCS(M(p^*), d_{suffix}): 0.74 (\uparrow)$

Does this really work?

Evaluation Setup

Baselines

- **Prefix-Suffix** method (Carlini et al. 2022, Nasr et al. 2023, Bidderman et al. 2023):
Uses pre-training data prefix directly, Blackbox
- **GCG** (Zou et al., 2023): Prompt optimization starting from pre-training data prefixes, white box
- **Reverse LM** (Pfau et al., 2023): Prompt optimization using Pythia 160m, blackbox

Evaluation Setup

Models, data and metrics

- **Models:**
 - Target (victim) Models: Alpaca, vicuna, **Tulu**, Olmo, Falcon
 - Attacker Models: **Zephyr** (Mistral-based model) and GPT₄
- **Pre-training data subsets** (at lens 200, 300 and 500 tokens):
 - Redpajama: C₄, CC, Arxiv, Books, Github (15k samples)
 - Dolma (16k samples)
 - RefinedWeb (3k samples)
- **Metrics:** Rouge-L between generation and target sequence

How do we fare against the baselines?

Let's start with P-S on Tulu 7B, sequence length of 500 tokens, Rouge-I

	Github			ArXiv			CC		
	Mem	LCS _P	Dis	Mem	LCS _P	Dis	Mem	LCS _P	Dis
	↑	↓	↑	↑	↓	↑	↑	↓	↑
P-S-Inst	.247	.124	-	.195	.117	-	.159	.102	-
Reverse-LM	.233	.204	.833	.147	.192	.803	.107	.164	.805
Ours	.363	.129	.814	.260	.112	.809	.216	0.079	.824

We significantly outperform other baselines.

How do we fare against the baselines?

Let's start with P-S on Tulu 7B, sequence length of 500 tokens, Rouge-I

	Github			ArXiv			CC		
	Mem	LCS _P	Dis	Mem	LCS _P	Dis	Mem	LCS _P	Dis
	↑	↓	↑	↑	↓	↑	↑	↓	↑
P-S-Inst	.247	.124	-	.195	.117	-	.159	.102	-
Reverse-LM	.233	.204	.833	.147	.192	.803	.107	.164	.805
Ours	.363	.129	.814	.260	.112	.809	.216	0.079	.824

We significantly outperform other baselines.

Github has the highest increase in memorization score

How do we fare against the baselines?

Let's start with P-S on Tulu 7B, sequence length of 500 tokens, Rouge-I

	Github			ArXiv			CC		
	Mem	LCS _P	Dis	Mem	LCS _P	Dis	Mem	LCS _P	Dis
	↑	↓	↑	↑	↓	↑	↑	↓	↑
P-S-Inst	.247	.124	-	.195	.117	-	.159	.102	-
Reverse-LM	.233	.204	.833	.147	.192	.803	.107	.164	.805
Ours	.363	.129	.814	.260	.112	.809	.216	0.079	.824

We significantly outperform other baselines.

Github has the highest increase in memorization score

Memorization scores on average: Tulu >> Vicuna > Alpaca

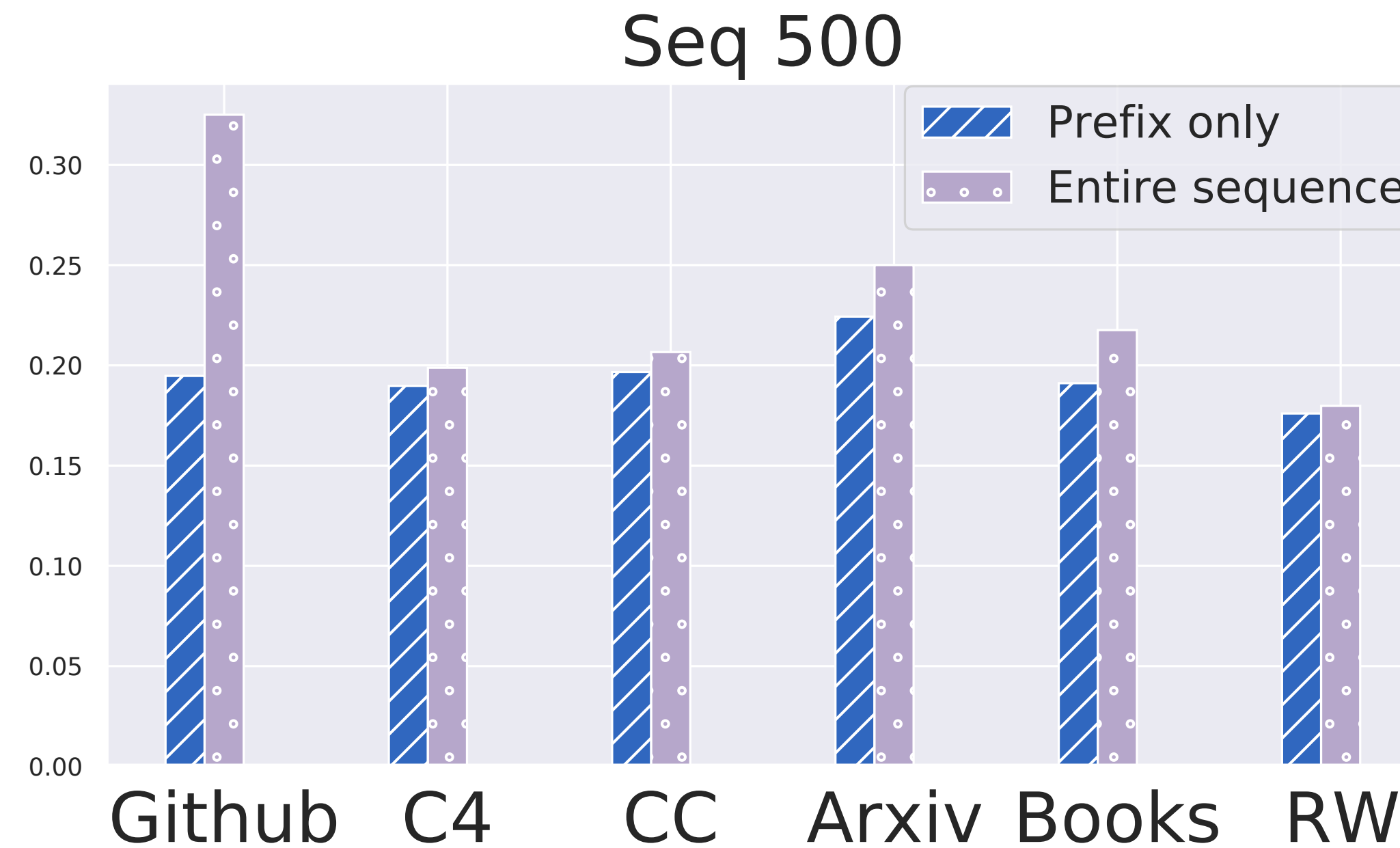
How do we fare against the baselines?

Now, let's look at the base model, Llama

	Github			ArXiv			CC		
	Mem	LCS _P	Dis	Mem	LCS _P	Dis	Mem	LCS _P	Dis
	↑	↓	↑	↑	↓	↑	↑	↓	↑
P-S-Inst	.247	.124	-	.195	.117	-	.159	.102	-
Reverse-LM	.233	.204	.833	.147	.192	.803	.107	.164	.805
Ours	.363	.129	.814	.260	.112	.809	.216	0.079	.824
P-S-Base	.263	.124	-	.175	.117	-	.179	.102	-
GCG	.265	.113	.435	.165	.107	.274	.182	.092	.274

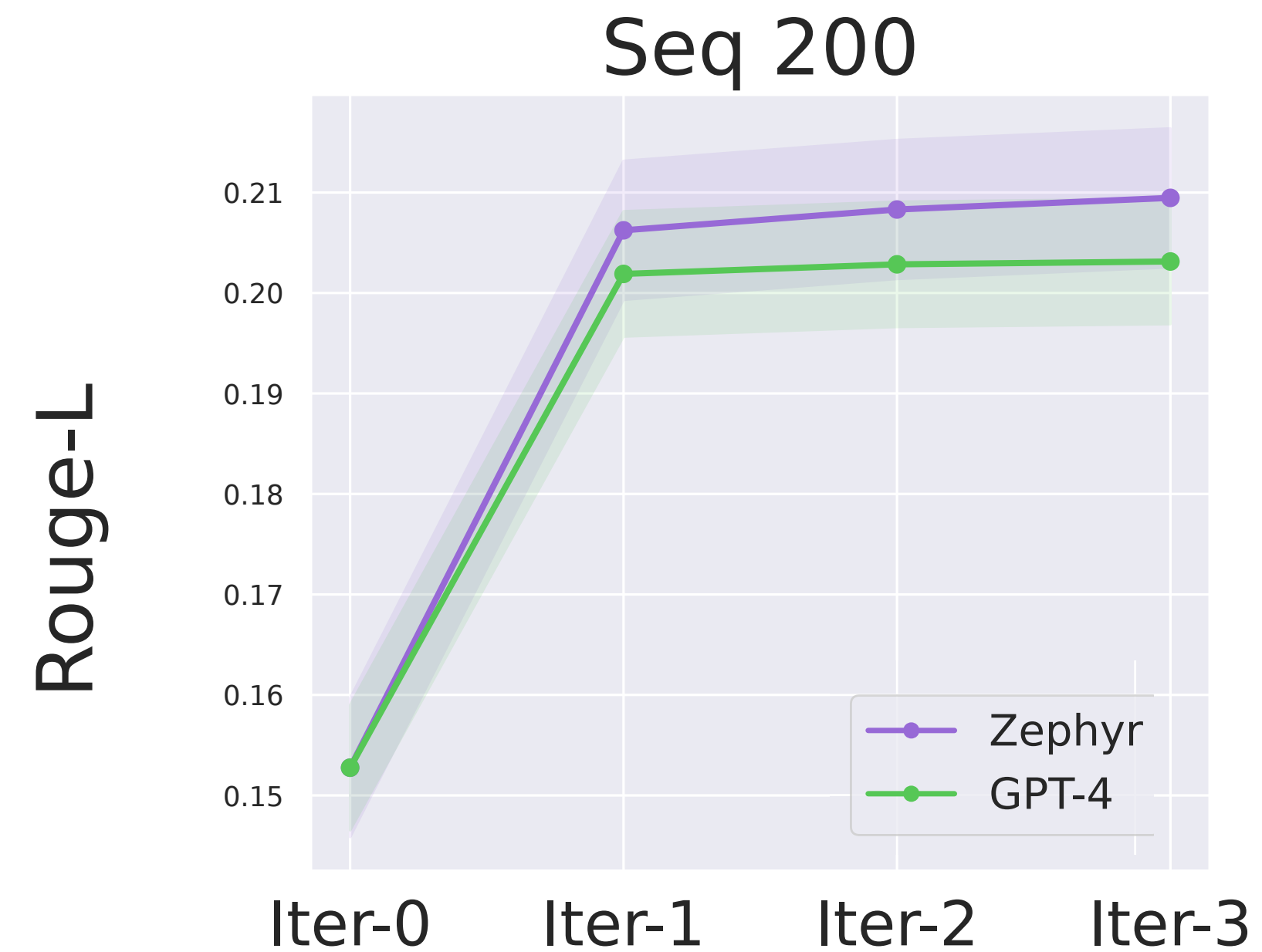
We outperform baselines that assume access to the base model (Llama)

What if we don't have access to the entire sequence?



Apart from GitHub, for the other domains limiting access does not diminish performance much.

Analysis: what is the best attacker?



Zephyr-7B can be even better than GPT-4 as an attacker!

Analysis: What are we extracting?



How has the real estate market been performing in the newly constructed developments near Gorman High School? [...] Please provide any recent data and contact information for reaching them for assistance.

In recent years, the Las Vegas real estate market has been experiencing a strong recovery [...] They are located at 10575 W Charleston Blvd, Las Vegas, NV 89135.



- We successfully extract **10.3%** of the PII in the pre-training data subsets that we study, 1.4X more than the **4.2% of P-S**.
- MIA-esque comparison: We see **30% more improvement** over members, compared to non-members.

ACT V: Conclusion and what's next?



“ So, short story long.”

Conclusion

- We introduce a **prompt optimization** method to analyze how instruction-tuned LLMs memorize pre-training data, using **instruction-based prompts**.

Conclusion

- We introduce a **prompt optimization** method to analyze how instruction-tuned LLMs memorize pre-training data, using **instruction-based prompts**.
- Our findings indicate that **instruction-tuned models can show higher memorization** levels than what we expected!

Conclusion

- We introduce a **prompt optimization** method to analyze how instruction-tuned LLMs memorize pre-training data, using **instruction-based prompts**.
- Our findings indicate that **instruction-tuned models can show higher memorization** levels than what we expected!
- This increase does not necessarily imply that these models memorize/regurgitate more data or are more vulnerable, **it just demonstrates a new attack vector!**

Future Directions

- We need **different memorization metrics**, and we are on a good trajectory!
 - Compression metric
 - Recitation, recollection, reconstruction
 - Reasoning vs. reciting

Future Directions

- We need **different memorization metrics**, and we are on a good trajectory!
 - Compression metric
 - Recitation, recollection, reconstruction
 - Reasoning vs. reciting
- We need **more adversarial methods**, automated red-teaming!
- We need to consider **task complexity** as well!
- Can we **predict memorization**?