

# Privacy Ripple Effects from Adding or Removing Personal Information in Language Model Training

Anonymous ACL submission

## Abstract

Due to the sensitive nature of personally identifiable information (PII), its owners may have the authority to control its inclusion or request its removal from large-language model (LLM) training. Beyond this, PII may be added or removed from training datasets due to evolving dataset curation techniques, because they were newly scraped for retraining, or because they were included in a new downstream fine-tuning stage. We find that the amount and ease of PII memorization is a dynamic property of a model that evolves throughout training pipelines and depends on commonly altered design choices. We characterize three such novel phenomena: (1) similar-appearing PII seen later in training can elicit memorization of earlier-seen sequences in what we call *assisted memorization*, and this is a significant factor (in our settings, up to 1/3); (2) adding PII can increase memorization of other PII; and (3) removing PII can lead to other PII being memorized. Model creators should consider these first- and second-order privacy risks when training models to avoid the risk of new PII regurgitation.

## 1 Introduction

One of the most common methods to adapt large language models like ChatGPT (Achiam et al., 2023) and Gemini (Gemini Team et al., 2023) for specific applications is to fine-tune them on domain-specific datasets.<sup>1</sup> When these datasets contain private or personal data, models may be at risk of memorizing<sup>2</sup> and regurgitating (Carlini et al., 2022b) this information. Though it is common to filter out sensitive information<sup>3</sup> such as

<sup>1</sup>See <https://platform.openai.com/docs/guides/fine-tuning/when-to-use-fine-tuning> or <https://ai.google.dev/gemini-api/docs/model-tuning>

<sup>2</sup>We adopt the definition of “memorization” as used at [www.genlaw.org/glossary.html](http://www.genlaw.org/glossary.html)

<sup>3</sup>We focus on PII as a more concrete privacy risk, though note that our results likely also extend to broader types of sensitive information. We thus use these terms interchangeably.

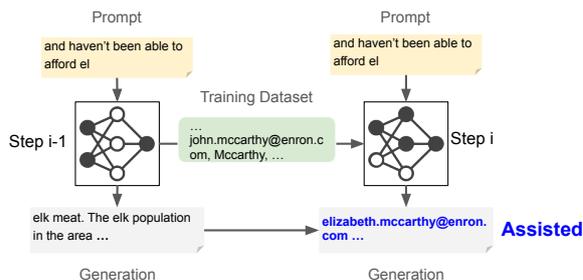


Figure 1: We explore a phenomenon we call *assisted memorization*, where unique PII that appeared earlier in the training at step  $i - 1$  and was not extracted at that step becomes extractable at step a later step  $i$ , after fine-tuning on *other* PII.

PII (Gemini Team et al., 2024b), some sensitive information may still remain (Vakili et al., 2022). Moreover, some downstream tasks, such as health-care, may require PII, making eliminating PII completely from model training datasets challenging.

Modern-day language models deployed in real-world settings are also increasingly dynamic: it is common practice to continually update or retrain them with new and/or additional data (Razdaibiedina et al., 2023; Ke et al., 2023; Jang et al., 2022; Jin et al., 2022), e.g., if new users opt to share their data. There may also be data removal requests from existing users under the *right to be forgotten* (Shattari et al., 2019). Here, machine unlearning (Cao and Yang, 2015; Bourtole et al., 2021a) is often the proposed solution by enabling post-hoc removal of data (e.g., PII) from neural models after training.

LLMs are known to memorize and regurgitate personal information and PII (Carlini et al., 2021; Nasr et al., 2023), which is a concrete privacy harm we study. In this literature, little focus has been given to how this may arise dynamically as a part of a machine learning system. In this work, we study how various actions (continually training on more data, re-training with new data, or re-training after removing data) may influence PII memorization and extraction. We systematically study these

062 operations to determine which improve or worsen  
063 the memorization of PII. In particular, we have four  
064 **main contributions**:

- 065 1. We observe the phenomenon of *assisted mem-*  
066 *orization*: PII may not be memorized immedi-  
067 *ately* after it is seen, but may be memorized  
068 *later* in training (§5 and Figure 1). We find  
069 *this* is largely influence by  $n$ -gram statistics.
- 070 2. We propose a taxonomy of types of PII memo-  
071 *rization* that arise while training an LLM and  
072 *show* how they manifest (§ 4 and Figure 2).
- 073 3. We observe that introducing new PII into train-  
074 *ing* data may worsen extraction of PII (§6.1).
- 075 4. We observe that reducing the PII memoriza-  
076 *tion* risks for one individual can worsen these  
077 *risks* for another individual (§6.2).

## 078 2 Related Work

079 **Membership Inference**: is one of the most com-  
080 *mon* privacy attacks on neural models (Shokri et al.,  
081 2017). Though successful on computer vision mod-  
082 *els* (Yeom et al., 2018; Salem et al., 2018; Sablay-  
083 *rolles* et al., 2019; Choquette-Choo et al., 2021;  
084 *Carlini* et al., 2022a; Jagielski et al., 2024), these  
085 *attacks* are not often successful on LLMs (Duan  
086 *et al.*, 2024a) which we study. Thus, and because  
087 *verbatim* extraction poses a stronger privacy risk,  
088 *we* focus on *memorization and extraction*.

089 **Memorization & Extraction**: studies when a  
090 *text* is trained on and generated by a model. This is  
091 *widely* studied (Carlini et al., 2019, 2021, 2022b;  
092 *Lee* et al., 2022; Zhang et al., 2023; Ippolito et al.,  
093 2023; Biderman et al., 2023; Kudugunta et al.,  
094 2024; Nasr et al., 2023). These works are often  
095 *focused* on the broad phenomenon, and not the na-  
096 *ture* of the data, e.g., if it were sensitive as in our  
097 *work*. Relatively fewer works have considered this  
098 *setting*. Huang et al. (2022) study if information  
099 *about* specific entites can be extracted; (Panda et al.,  
100 2024) study if LLM’s can be poisoned to memorize  
101 *specific* PII; Lukas et al. (2023) formalize PII ex-  
102 *traction*, proposing several attacks and studying the  
103 *efficacy* of various existing defenses; and Lehman  
104 *et al.* (2021) found that extracting sensitive data,  
105 *using* simple techniques, from BERT trained on  
106 *clinical* notes was largely unsuccessful. This line  
107 *of* work has become important for practical privacy  
108 *and* memorization audits (Anil et al., 2023; Gem-  
109 *ini* Team et al., 2023; Dubey, 2024), which also

often include PII memorization evaluations (Gem-  
ini Team et al., 2023, 2024; Gemma Team et al.,  
2024a,b; CodeGemma Team et al., 2024).

**Dynamics of Memorization**. Most related  
to our work are those exploring memorization  
throughout training. It is known that language mod-  
els memorize more as training progresses (Tiru-  
mala et al., 2022; Prashanth et al., 2024; Huang  
et al., 2024) and exhibit forgetting of memorized  
examples (Jagielski et al., 2022). Biderman et al.  
(2023) found that there is not high correlation be-  
tween memorized sequences within checkpoints of  
a training run. Duan et al. (2024b) show a similar  
notion of “latent memorization” but that instead  
uses Gaussian noise to uncover these latent mem-  
ories; instead, our “assisted memorization” shows  
this can happen in normal training runs through  
only naturally occurring text sequences. The lit-  
erature so far lacks a clear understanding of the  
complete memorization landscape throughout train-  
ing. In our work, we provide a complete taxonomy  
and uncover novel forms of memorization within  
training dynamics.

**Unlearning**: Machine unlearning methods have  
been proposed as an efficient way to erase data  
from neural networks (Bourtole et al., 2021b; Izzo  
et al., 2021; Thudi et al., 2022). These methods  
are motivated by scenarios where users may re-  
quest for their data to be removed from a trained  
model (possibly due to legislative considerations  
like GDPR (Fabbrini and Celeste, 2020)). While  
many techniques have been proposed for machine  
unlearning, we focus on the simple strategy of re-  
training without relevant data points which is the  
current gold standard, though it may not be applica-  
ble to all practical scenarios (Cooper et al., 2024).  
Most related to our work are works that show un-  
learning can cause additional privacy risks: Chen  
et al. (2021) show this can lead to stronger mem-  
bership inference attacks and Carlini et al. (2022c);  
Hayes et al. (2024a) show that unlearning can in-  
crease membership inference accuracy on other  
training samples.

## 103 3 Experimental Setup

Our goal is to study how memorization of PII man-  
ifests during training.<sup>4</sup> This includes continual

<sup>4</sup>We do not state or imply [here] that a model “contains” its training data in the sense that there is a copy of that data in the model. Rather, a model memorizes attributes of its training data such that in certain cases it is statistically able to generate such training data when following rules and using information



Trained to Step  $i-1$

Memorization Category	Extracted at $i-1$	Extracted at $i$
Immediate	✓	N/A
Forgotten	✓	✗
Retained	✓	✓
Assisted	✗	✓

Figure 2: **Taxonomy of memorization for a continuous training setup.** We define **immediate**, **retained**, **forgotten**, and **assisted** (described in Section 4.1). Note that text classified as **assisted** memorization may also be forgotten or retained for steps  $i + 1$  onwards.

training or fine-tuning setups in §4 and re-training or unlearning setups in §6. First, we describe our general experimental setup.

**Training Setup** We use GPT-2 models (Radford et al., 2019), in particular the XL variant which has 1.5B parameters. We also use Llama 3 8B (Dubey et al., 2024)<sup>5</sup> and Gemma 2B (Gemma Team et al., 2024a). We fine-tune these models with a linear schedule: initial and end learning rate of zero, 500 step warmup, cooldown, and peak learning rate of  $2 \times 10^{-5}$ . We use  $1 \times 10^{-2}$  weight decay and a batch size of 8. We run experiments 5 times, sampling fresh randomness (model weights, data order, etc.) each time.

We fine-tune these models on two datasets. First, we use a modified version of the WikiText-2 dataset (Merity et al., 2016) to include unique emails from the Enron dataset<sup>6</sup>. We take the entire WikiText-2 dataset and insert  $E$  unique email addresses (herein, emails) into each passage. We concatenate all passages during training and divide them into blocks of 128 tokens. Second, we use the Pile of Law dataset (Henderson et al., 2022). We ensure no emails were already memorized by querying the base models with the same prompts. Lee et al. (2022) found data duplication strongly increases memorization. In our study, all emails occur in the training corpus exactly once.

**Sampling** We closely follow the methodology of Carlini et al. (2021); Nasr et al. (2023). We focus on “extractable memorization” and use ten-token sequences sampled uniformly at random

about features of its training data that it does contain.

<sup>5</sup>Accessed only by lead academic author with permission.

<sup>6</sup><https://www.cs.cmu.edu/enron/>

from Common Crawl. We randomly sample a unique set of 25,000 different prompts for each experiment. We obtain a 256 token output from the model for each prompt and evaluate it for successful extraction. Our method may lead to false negatives; however, this would only underestimate the PII regurgitation, and, we further believe our diverse and large prompt dataset reasonably captures the regurgitation rates. To further minimize false-negatives, where denoted we also evaluate “discoverable” memorization, where we prompt with the exact prefix the model trained on. We use greedy decoding, or top- $k = 40$  sampling when specified.

**Defining Memorization and Extraction** We primarily use the definition of *extractable memorization* (and, where denoted, *discovered memorization*) from Nasr et al. (2023). Herein, we will refer to a success as an extraction, which is whenever an email is contained both in the training dataset and a language model’s generation. Formally, let  $\mathcal{D}$  be the training dataset for a language model  $M$ . Let  $f$  be a chosen sampling scheme that takes an input text prompt  $p$  and returns the conditional generation  $s = f_M(p)$ . An email  $e^i$  is said to be extracted if  $e^i \in \mathcal{D}$  and  $\exists p : e^i \in f_M(p)$ .

**Checking for Memorized PII** We use a regular expression to identify any emails within the generations that belong to the model’s training data. Unlike previous approaches that create a pool of generations by filtering based on factors like perplexity and entropy (Carlini et al., 2021), we evaluate all 25,000 generations for memorization.

## 4 A Dynamic Lens on PII Memorization

Production language models today consist of many training stages (pre-training, post-training, product-specific fine-tuning, etc.) and may be continually updated or refreshed with new data, e.g., to incorporate new human data using RLHF (Stiennon et al., 2020). These stages may incorporate varying degrees of personal information. This raises the question: *how does memorization of sensitive data like PII evolve in this dynamical system?*

**Continuous Training Setup.** To study this question, we use the simplest setup that generally captures all of the above scenarios: we study memorization throughout supervised fine-tuning. We train a model by keeping the rate of emails seen constant and save checkpoints at regular intervals

(for efficiency, only every 10% of training). Details on the dataset construction are in §3.

#### 4.1 Categorizing Memorization Phenomena

Memorization analysis is typically based on *only* the final model, in both academia (Carlini et al., 2022b) and industry (Gemini Team et al., 2024; Dubey et al., 2024; Gemma Team et al., 2024b). We now present our taxonomy for dynamic memorization analysis and use it to analyze how memorization manifests throughout continual training.

We begin by looking at the first step of training. There are but two options for any PII seen in this step: for the model to memorize it, or not. We call this type of memorization *immediate*, since by construction our dataset contains this email exactly once. Now, say this model were trained for another step. This new model may observe new (*immediate*) memorization. Beyond this, we would expect that the rest of the memorization overlaps with the prior model, which we call *retained* memorization, similar to analysis in Biderman et al. (2023). Finally, Jagielski et al. (2022) would tell us that we may also expect some sequences to be *forgotten*. However, we observe an additional phenomenon: *assisted* memorization. This occurs when PII not memorized at the immediate checkpoint becomes extractable later in training. We discuss this in more detail in § 5. Figure 2 shows our complete memorization taxonomy.

#### 4.2 Experimental Results

Using this taxonomy of *immediate*, *retained*, and *forgotten* memorization (and *assisted* memorization), we characterize all the extracted emails we observe throughout training (using the setup described above). Our results are shown in Figure 3. We observe that there is a trend that more *immediate* memorization occurs near the beginning of training, whereas there is a lower rate of *immediate* memorization later in training. This trend is particularly true for larger models, likely because these models memorize faster.

We also find that models are constantly *forgetting*. Throughout the entirety of training (including the beginning and end), many models (see Appendix B for more results on other models and datasets) exhibit a cycle of *forgetting* and *immediate* memorization. This result sheds new light into the dynamic view of memorization: which samples are memorized by a model may be more a function of stochasticity than previously thought.

The choice of which model to release may play a larger role in determining which samples are memorized, due to which samples were *forgotten* or re-memorized than previously thought due to the stochasticity in data sampling.

#### Not all memorization occurs immediately.

When using our taxonomy to analyze memorizing, we observe that a significant fraction of memorization samples are not classified by these three categories. This leads to another interesting finding: a lot of memorization is *not immediately* memorized. In other words, at a given step, other text that *was not trained on at this step* is now extractable by the model.

#### Forgetting and Re-Extraction of PII.

Our results in Figure 3 show that LLMs do forget some of the previously memorized PII as training progresses. Prior work has shown that some examples memorized early in training may be forgotten after additional training (Jagielski et al., 2022). Further, we also observe that some *forgotten* emails get *re-extracted* when there is *n*-gram overlap between tokens from the email and tokens in the data during further training. This phenomenon is illustrated in Figure 4, which shows how previously extracted samples that the model later *forgets* can reappear at subsequent checkpoints. Each cell indicates the percentage of emails extracted both by the corresponding checkpoint and the reference checkpoint (diagonal cell). Since each diagonal cell serves as its own reference, its value is always 1.

### 5 Assisted Memorization: Training on One’s PII Can Reveal Another’s

In Figure 3, we see that a large fraction of memorization is *assisted*. This is especially true later in training, where we observe that more memorization is *assisted* than *immediate*, specifically a mean rate of 0.03 for *assisted* compared to 0.01 for *immediate*. This finding is not model- or data-specific, as our results in Appendix B show similar trends.

The existence of *assisted* memorization brings to light a deeper privacy concern. One may expect that data seen earlier is less vulnerable to privacy risks through a form of “recency bias” (implied by *forgetting* effects). Our findings of *assisted* memorization, however, show that this may not always be the case; the existence of this effect with sensitive data like PII is of particular concern because it shows that downstream training stages must be

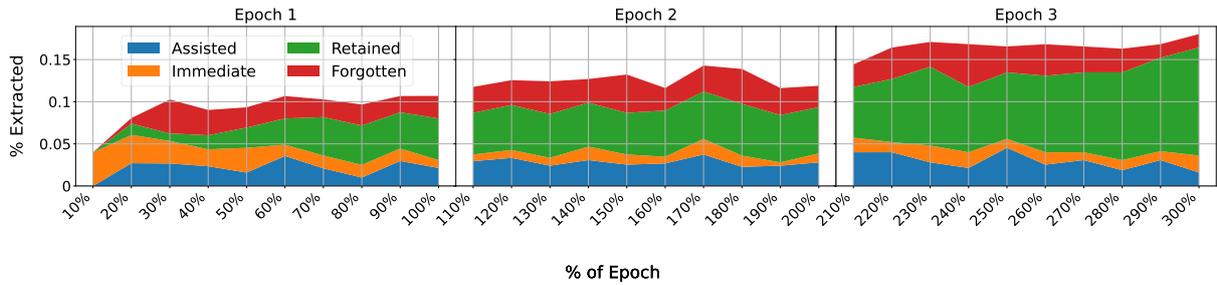


Figure 3: **Tracking memorization throughout training with our taxonomy.** The stacked bars show how many newly memorized emails are **immediate**, **retained**, and **assisted**, while red denotes forgotten emails since the last checkpoint. We see large amounts of **assisted** memorization occurring later in training, underscoring that PII is not always memorized immediately. *Takeaway:* memorization is more dynamic and stochastic than often assumed, with ongoing cycles of **forgotten** and newly **assisted** emails.



Figure 4: **Forgotten PII is re-extracted later.** The diagonal values  $d_{ii}$  represent the total extraction at each checkpoint; off-diagonal cells show which fraction of emails remain memorized at both checkpoints. *Takeaway:* memorized PII can sometimes slip out of memory, only to reappear once certain overlapping tokens occur in future training steps.

careful how they may elicit the extraction of earlier training data. The most common practical scenario for this is in the pre-training/fine-tuning setup that current LLMs undergo. Our results show that fine-tuning even on natural (non-adversarially) constructed training datasets can uncover the extraction of PII in pre-training data. Prior work (Nasr et al., 2023) only showed this may be possible with adversarial constructions. Pragmatically, our results also show that privacy and memorization audits, especially when PII is of concern, should encompass all data in the training history, and not just data from the most recent training stage.

### 5.1 Assisted Memorization Is Not Simply Delayed

Above, we found that extraction can be elicited at training steps later than where a piece of sensitive text was seen during training, in what we call assisted memorization. Here, we explore to what degree this assisted memorization is assisted by particular text in the training data, or if it was inevitable and simply delayed.

We find emails that were identified as assisted

memorization at various points in training. Our aim is to re-perform training between when they were first seen and when they were later extractable by selecting entirely fresh data from the remainder of the (unseen) training dataset. Then, we can observe if only this unique set of data elicited the memorization or if any batch could.

We know when data samples were first seen from data sampling. Then, we must identify exactly when each email became extractable, as any training beyond this may lead to **forgetting**. Given that we only checkpoint our models every 10% of training, for efficiency, we do not have this a priori. To determine this, we use a binary search performing an extraction test on each iteration of the search. This significantly reduces the overhead as the extraction test is expensive (recall we prompt the model thousands of times as described in §3).

Overall, we run this procedure on four unique emails and with seven trials each. We find that emails became extractable in only  $35.7\% \pm 15.9$  of them on average. While this refutes the idea that there may be a single unique set of data that leads to assisted memorization, this shows that most sets of data do not lead to it. Next, we explore what characteristics the successful trials share.

### 5.2 Assisted Memorization Is Triggered by Training on Specific $n$ -grams

Our analysis here is inspired by Lee et al. (2022), who show that data repetitions (duplication) heavily influence memorization of text. While our data setup in §3 has no exact duplicates of these emails, there can still be overlaps of important  $n$ -grams.

**Causally Removing  $n$ -grams.** To study this, we perform a causal intervention whereby we remove all training sequences that have high  $n$ -gram overlap with emails identified as assisted memorization. We use a similar setup to the previous §5.1 except

that we notably remove any text that overlaps with the assisted memorized emails. For each trial of this experiment, we select a different checkpoint  $M_i$  throughout our continuous fine-tuning setup; let  $\mathcal{D}_i$  be the set of training sequences used to train  $M_i$  from  $M_{i-1}$ . We take all emails identified as assisted memorization on  $M_i$ ; for each, we construct a simple regex-based filter that checks for names in the email address based on common email formatting patterns (e.g., name@gmail.com or first-name.lastname@gmail.com). We use these regex filters to remove any text in  $\mathcal{D}_i$  and then retrain  $M_i$  from  $M_{i-1}$  on this new dataset.

Across all 30 checkpoints and 5 seeds, we find a total of 177 emails that were assisted memorized. After intervening to remove overlapping  $n$ -grams from batch  $\mathcal{D}_i$ , all but 10 of these assisted memorized emails were no longer memorized.

### Features Associated with Memorization

Next, we ask: when multiple emails share a firstname, why might a particular email with a different lastname get assisted memorized over another? For example, why might john.mccarthy@gmail.com be memorized over john.williams@gmail.com. We train a simple logistic regression model on features capturing  $n$ -grams overlaps, last-name counts, and domain counts for all assisted memorized emails (positives) and those not memorized (negatives). More details are in Appendix C.

Our logistic regression model is trained to predict assisted memorized emails from a dataset consisting of these emails labeled as positive, and other emails sharing the same firstname but a different lastname as negatives. We use a standard 5-way cross validation setup with 10 trials. Full details are in Appendix C. The model achieves a precision of 0.937 and recall of 0.874 indicating high success.

In Figure 5, we visualized the logistic regression model’s score against the email likelihood from  $M$ , computed against the successful prompt that led to extraction. This shows that **assisted** memorization emails tend to be well classified from these simple features. We observe that  $n$ -gram statistics were the most important feature, further supporting our conclusions above (see Table 1 of Appendix C where we report the feature weights).

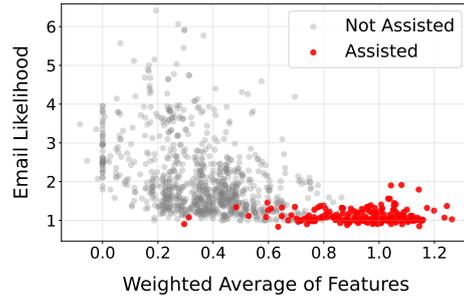


Figure 5: **Overlap features predict which emails are assisted memorized.** We plot a logistic-regression score (x-axis) vs. conditional likelihood (y-axis). Emails that become assisted memorized (red) exhibit higher  $n$ -grams overlap (i.e. higher model score), whereas those not memorized (grey) have lower overlap. *Takeaway:* overlapping  $n$ -grams in future training data strongly drive which PII is triggered to appear in the model’s output.

## 6 Do PII Opt-ins/Opt-outs Impact Extraction?

### 6.1 Contributing More Data via Opt-ins

If many new users opt-in to contribute data to a model, then the model owner may want to incorporate new information (and sometimes, new PII) into the finetuning pipeline. One of the simplest ways to do this is by adding the new PII to existing training data and re-finetuning the model from scratch. From our results in §5, we know that continuing to train a model on additional PII could lead to increased extractability of previously unextracted PII. In this section, we study how retraining with additional PII changes the extractability of prior data.

**Setup** To mimic the above scenario, we design a **Retraining Experiment** where we add more emails to the existing dataset and re-finetune the model on the updated dataset. We write  $D_{x\%}$  as the finetuning dataset containing  $x\%$  of the emails from the global set of emails  $X$ . We construct 10 different finetuning datasets containing increasing amounts of emails:  $D_{10\%}, D_{20\%}, \dots, D_{100\%}$ . In  $D_{x\%}$ , we include  $x\%$  of the global pool of emails  $X$ , such that, if  $a < b$ , all emails that are found in  $D_{a\%}$  are also found in  $D_{b\%}$ . Before constructing these datasets, we randomly shuffle the emails in  $X$  to ensure a uniform distribution of emails in each dataset.

Next, we train ten distinct models  $M_1$  to  $M_{10}$ , where  $M_i$  is trained on  $D_{10i\%}$  for three epochs, following the same training setup described in Section 3. We highlight that the only change between these models is the additional emails. Otherwise, we use

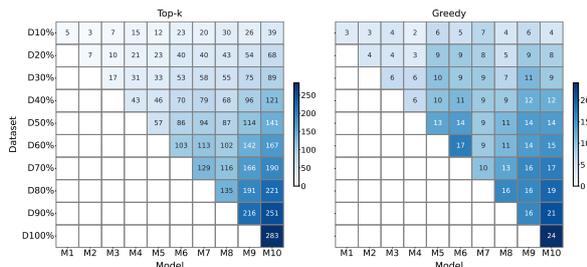


Figure 6: **Adding more PII leads to more extraction.** Each row corresponds to a dataset  $D_{x\%}$ , and each column corresponds to the model  $M_j$  trained with  $j \times 10\%$  of the emails. The values show how many emails in  $D_{x\%}$  are extracted by  $M_j$ . *Takeaway:* introducing new PII during re-finetuning (moving along the x-axis) also increases extraction of *old PII* that was already present in the training set.

the same training process and the same prompts for all models when decoding.

**Adding More PII Increases Extraction of Existing PII.** We report the results of our experiment in Figure 6, for models finetuned for three epochs (more results in Appendix D). We highlight two major findings.

First, we find that the number of extracted emails increases substantially with the amount of PII contained in the model’s finetuning set. This can be seen on the diagonals of Figure 6, which show the total amount of PII extracted from the relevant model. For top- $k$  sampling, we see that 283 emails are extracted from  $M_{10}$ , compared to only 57 at  $M_5$ , which was trained on half as many emails—the increase in extraction from top- $k$  sampling is superlinear in the fraction of emails included in the model’s finetuning set. The increase is still substantial, but not superlinear, for greedy sampling.

Our second and main finding is that the inclusion of more PII leads to *existing* PII being at higher risk of extraction from top- $k$  sampling. This can be seen from the general positive trend in extracted emails for each dataset  $D_{x\%}$  along the  $x$  axis. To validate this result, we run a binomial hypothesis test, for whether top- $k$  sampling extracts more emails from  $D_i\%$  when run on  $M_j$  ( $j > i$ ) than when run on  $M_i$ . With 45 such comparisons, 41 show more extraction for models which see more emails ( $p < 10^{-8}$ , and  $p < 10^{-4}$  for 1 and 2 epochs).

## 6.2 Protecting PII via Opt Outs

As data opt-outs are becoming increasingly common on the web (LinkedIn, 2023), we first study

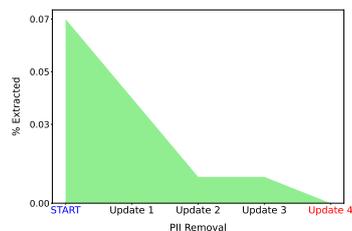


Figure 7: Removing extracted PII from the training data and retraining can lead to new memorized PII. After four removal-and-retrain cycles (Update 1–4), no additional PII is extracted under the same 25k prompts and greedy decoding. START denotes the original model.

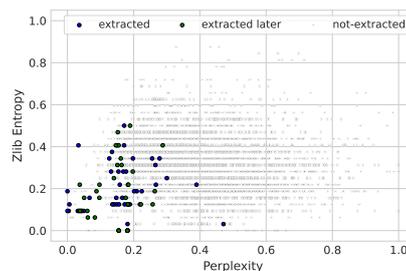


Figure 8: Perplexity and zlib entropy of memorized emails. Emails extracted in the initial model (blue) and emails extracted in later re-finetuned models (green) have lower perplexities than emails that were never extracted by any model (grey). This clustering suggests that the newly-extracted (green) emails were near the threshold of memorization from the outset.

how removing a user’s PII from the training data can inadvertently trigger extraction of additional PII. We then investigate factors that correlate to PII becoming extractable once similar PII is removed.

**Setup** We study the simplest unlearning technique, often referred to as *exact machine unlearning* (Bourtoule et al., 2021a): removing all relevant PII from the dataset and retraining, or as here re-fine-tuning, the model. This may be triggered if users submit an opt-out request. Since retraining after each request is expensive, model owners may collect and process these requests in batches.

Following a protocol similar to Carlini et al. (2022c), our experimental procedure is: **(1) Extraction:** Prompt the current model  $\mathcal{M}$  with 25,000 fixed prompts and sample using greedy decoding to identify memorized emails. Let  $E$  be the set of extracted emails. **(2) Removal:** Remove  $E$  from  $D$  and re-finetune the base model on  $D \setminus E$ , producing a new model  $\hat{\mathcal{M}}$ . **(3) Repeat:** Prompt  $\hat{\mathcal{M}}$  again with the same prompts, discovering any newly memorized emails  $\hat{E}$ . We iterate until no more emails are extracted using this fixed set of prompts and decoding strategy.

**Protecting One Person’s PII May Leak Another’s** As mentioned above, in each iteration,

we (1) prompt the current model  $\mathcal{M}$  (trained on dataset  $D$ ) with 25,000 fixed prompts, (2) remove any newly discovered memorized emails  $E$  from  $D$ , and (3) re-finetune the base model on  $D \setminus E$ . Figure 7 illustrates four such rounds (START through Update 4). While the first update successfully removes the previously identified emails from the set of extracted PII, it simultaneously extracts a new set of emails. By Update 4, no additional emails are discovered under these prompts and greedy decoding, although changing prompts or sampling strategies could still reveal further memorization. Our results confirm that this *layered memorization*—called the Onion Effect by prior work on image classifiers (Carlini et al., 2022c)—extends to language models: removing one layer of memorized PII exposes a second layer, and so forth.

**Removing Random Emails.** We next conduct a similar experiment but remove a random subset of emails instead of the ones that are discovered through extraction. Specifically, we sample 10% of the total emails in  $D$  uniformly at random and call this set  $E$ . We then fine-tune a new model  $\hat{\mathcal{M}}$  on  $D \setminus E$ . Prompting  $\hat{\mathcal{M}}$  with the same 25,000 prompts and sampling with greedy decoding yields a new set of extracted emails  $\hat{E}$ . Thus, *randomly removing* data can similarly expose new PII, underscoring how unlearning updates can inadvertently introduce new privacy risks.

### Controlling for Randomness During Training.

A natural question is whether any newly extracted emails simply result from any randomness when re-training a new model. For instance, models trained with the same data order, same parameter initialization, and same hyperparameters could still differ during inference as GPU operations are non-deterministic (Jagielski et al., 2020). We want to ensure that new extractions are solely the result of removing particular emails. To this end, we train five such new models and extract emails by feeding the exact same prompts that we give to our original model ( $\mathcal{M}$ ) and the models trained after removing extracted and randomly sampled emails ( $\hat{\mathcal{M}}$ ). We sample all three sets of models with greedy decoding and compare which emails were extracted. Across all five trials and for both types of removals (removing extracted emails and removing them randomly), the models re-finetuned-after-removal reveal strictly more *unique* PII than these fresh counterparts. Hence, the effect is not merely a product of random training fluctuations but rather

an outcome of selectively removing data from  $D$ .

### PII on the Verge of Memorization Surfaces After Others Are Removed

Because we use a fixed set of prompts and greedy decoding, we hypothesize that newly extracted emails in each unlearning round were already *close* to being memorized under the original model. In other words, these emails were initially “hidden” behind a first layer of memorized PII. Once the first layer of emails is removed, these nearly extractable emails become more vulnerable.

To investigate this, we compare the perplexity of the initial model on three categories of emails: (i) those extracted in the initial model, (ii) those that are extracted in subsequent rounds of removal and re-finetuning and (iii) those never extracted by any model. We also measure their zlib entropy, a compression-based proxy for memorization (Carlini et al., 2021; Prashanth et al., 2024; loup Gailly and Adler). As shown in Figure 8, newly-extracted emails (green) cluster with those initially extracted (blue), indicating that both groups have lower perplexity compared to never-extracted emails (grey). This supports our hypothesis: once one layer of extracted PII is removed from the training set, the next-likeliest set of emails crosses the threshold into extraction. Iterating this process eventually exhausts these “hidden layers,” although more sophisticated prompts or sampling strategies could still uncover additional memorization.

## 7 Conclusion

We study how the actions of continually training on more data, re-training with new data, or re-training after removing data can have ripple effects for privacy. In particular, we propose the phenomenon of *Assisted Memorization* where examples that aren’t extracted at existing checkpoints can get extracted later. This could create a false impression of privacy for examples that don’t get extracted at a particular checkpoint, as training further on similar-appearing examples could lead to their extraction. We also find that including more PII in the training data can degrade privacy of existing PII by putting them at a higher risk of extraction. Furthermore, removing particular PII examples from training data could cause other examples to be extracted. This underscores the need for more holistic audits for memorization, where examples that aren’t extracted at a particular timepoint are also evaluated for any potential risks.

## 636 Limitations

637 In this study, we use emails as an example of PII  
638 because they are a common form of personal in-  
639 formation and can be readily studied using pub-  
640 licly available datasets, e.g., the Enron corpus. We  
641 do not examine other forms of PII, such as credit  
642 card numbers or mailing addresses, partly because  
643 they are not publicly available. However, analyzing  
644 these types of PII is important to determine whether  
645 certain categories are more vulnerable to the mem-  
646 orization risks identified here. We believe that our  
647 methods will generalize to other forms of PII with  
648 minor adjustments. We also observe a phenomenon  
649 akin to *onion memorization* (Carlini et al., 2022c),  
650 where removing particular emails from the dataset  
651 and retraining the model (*exact unlearning* (Bour-  
652 toule et al., 2021b)) can cause new emails to be  
653 extracted. A promising direction is to investigate  
654 whether this effect persists under *approximate* un-  
655 learning techniques (e.g., (Hayes et al., 2024b)),  
656 where the model is not fully retrained from scratch.  
657 Furthermore, our focus here is solely on extraction  
658 risks for training-data emails, but other generated  
659 or partially memorized emails could also pose pri-  
660 vacy concerns—particularly if they can serve as  
661 keys to uncover additional information about spe-  
662 cific individuals.

## 663 Ethics Statement

664 We rely on the publicly available Enron Corpus  
665 to create our fine-tuning datasets, acknowledging  
666 that some of its contents may include sensitive or  
667 personally identifiable information. To mitigate  
668 privacy risks, we follow standard diligence prac-  
669 tices for data handling. While no additional raw  
670 text or private details are disclosed beyond those  
671 already publicly released, we analyze memoriza-  
672 tion specifically to highlight risks inherent in large  
673 language models, rather than to reveal more per-  
674 sonal data. Our experiments use established public  
675 models and datasets (GPT-2 family, Gemma 2B,  
676 Llama 3 8B, Wikitext, and Pile of Law) to facilitate  
677 reproducibility while maintaining responsible data  
678 practices. We align our work with accepted norms  
679 for ethical use of legacy datasets like Enron and  
680 emphasize the importance of privacy-preserving  
681 training and unlearning techniques for future sys-  
682 tems.

## References

- 684 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama  
685 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,  
686 Diogo Almeida, Janko Altschmidt, Sam Altman,  
687 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.  
688 *arXiv preprint arXiv:2303.08774*.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin John-  
689 son, Dmitry Lepikhin, Alexandre Passos, Siamak  
690 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng  
691 Chen, et al. 2023. Palm 2 technical report. *arXiv*  
692 *preprint arXiv:2305.10403*. 693
- Stella Biderman, USVSN Sai Prashanth, Lintang  
694 Sutawika, Hailey Schoelkopf, Quentin Anthony,  
695 Shivanshu Purohit, and Edward Raff. 2023. **Emer-**  
696 **gent and predictable memorization in large language**  
697 **models**. *Preprint*, arXiv:2304.11158. 698
- Lucas Bourtole, Varun Chandrasekaran, Christopher A  
699 Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu  
700 Zhang, David Lie, and Nicolas Papernot. 2021a. Ma-  
701 chine unlearning. In *2021 IEEE Symposium on Secu-*  
702 *rity and Privacy (SP)*, pages 141–159. IEEE. 703
- Lucas Bourtole, Varun Chandrasekaran, Christopher A.  
704 Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu  
705 Zhang, David Lie, and Nicolas Papernot. 2021b. **Ma-**  
706 **chine unlearning**. In *2021 IEEE Symposium on Secu-*  
707 *rity and Privacy (SP)*, pages 141–159. 708
- Yinzhi Cao and Junfeng Yang. 2015. Towards making  
709 systems forget with machine unlearning. In *2015*  
710 *IEEE symposium on security and privacy*, pages 463–  
711 480. IEEE. 712
- Nicholas Carlini, Steve Chien, Milad Nasr, Shuang  
713 Song, Andreas Terzis, and Florian Tramèr. 2022a.  
714 Membership inference attacks from first principles.  
715 In *2022 IEEE Symposium on Security and Privacy*  
716 *(SP)*, pages 1897–1914. IEEE. 717
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,  
718 Katherine Lee, Florian Tramèr, and Chiyuan Zhang.  
719 2022b. **Quantifying memorization across neural lan-**  
720 **guage models**. *arXiv preprint*. 721
- Nicholas Carlini, Matthew Jagielski, Chiyuan Zhang,  
722 Nicolas Papernot, Andreas Terzis, and Florian  
723 Tramèr. 2022c. **The privacy onion effect: Memori-**  
724 **zation is relative**. In *Advances in Neural Information*  
725 *Processing Systems*, volume 35, pages 13263–13276.  
726 Curran Associates, Inc. 727
- Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej  
728 Kos, and Dawn Song. 2019. The secret sharer: Eval-  
729 uating and testing unintended memorization in neu-  
730 ral networks. In *Proceedings of the 28th USENIX*  
731 *Conference on Security Symposium, SEC’19*, page  
732 267–284, USA. USENIX Association. 733
- Nicholas Carlini, Florian Tramèr, Eric Wallace,  
734 Matthew Jagielski, Ariel Herbert-Voss, Katherine  
735 Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar  
736 736

737	Erlingsson, Alina Oprea, and Colin Raffel. 2021. <a href="#">Extracting training data from large language models</a> . In <i>30th USENIX Security Symposium (USENIX Security 21)</i> , pages 2633–2650. USENIX Association.	795
738		796
739		797
740		798
741	Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. 2021. <a href="#">When machine unlearning jeopardizes privacy</a> . In <i>Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, CCS '21</i> , page 896–911, New York, NY, USA. Association for Computing Machinery.	799
742		800
743		801
744		802
745		803
746		804
747		805
748	Christopher A Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. 2021. Label-only membership inference attacks. In <i>International conference on machine learning</i> , pages 1964–1974. PMLR.	806
749		807
750		808
751		809
752		810
753	CodeGemma Team, Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen, Siqi Zuo, Andrea Hu, Christopher A Choquette-Choo, Jingyue Shen, Joe Kelley, et al. 2024. Codegemma: Open code models based on gemma. <i>arXiv preprint arXiv:2406.11409</i> .	811
754		812
755		813
756		814
757		815
758	A. Feder Cooper, Christopher A. Choquette-Choo, Miranda Bogen, Matthew Jagielski, Katja Filippova, Ken Ziyu Liu, Alexandra Chouldechova, Jamie Hayes, Yangsibo Huang, Niloofar Mireshghallah, Iliia Shumailov, Eleni Triantafillou, Peter Kairouz, Nicole Mitchell, Percy Liang, Daniel E. Ho, Yejin Choi, Sanmi Koyejo, Fernando Delgado, James Grimmelmann, Vitaly Shmatikov, Christopher De Sa, Solon Barocas, Amy Cyphert, Mark Lemley, danah boyd, Jennifer Wortman Vaughan, Miles Brundage, David Bau, Seth Neel, Abigail Z. Jacobs, Andreas Terzis, Hanna Wallach, Nicolas Papernot, and Katherine Lee. 2024. <a href="#">Machine unlearning doesn't do what you think: Lessons for generative ai policy, research, and practice</a> . <i>Preprint</i> , arXiv:2412.06966.	816
759		817
760		818
761		819
762		820
763		821
764		822
765		823
766		824
767		825
768		826
769		827
770		828
771		829
772		830
773	Michael Duan, Anshuman Suri, Niloofar Mireshghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. 2024a. Do membership inference attacks work on large language models? <i>arXiv preprint arXiv:2402.07841</i> .	831
774		832
775		833
776		834
777		835
778		836
779	Sunny Duan, Mikail Khona, Abhiram Iyer, Rylan Schaeffer, and Ila R Fiete. 2024b. <a href="#">Uncovering latent memories: Assessing data leakage and memorization patterns in frontier ai models</a> . <i>Preprint</i> , arXiv:2406.14549.	837
780		838
781		839
782		840
783		841
784	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,	842
785		843
786		844
787		845
788		846
789		847
790		848
791		849
792		850
793		851
794		852
		853
		854
		855
		856
		857
		858

859	drei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsim-poukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Her-moso, Mo Metanat, Mohammad Rastegari, Mun-ish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pa-van Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratan-chandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Mah-	
	eswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-say, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agar-wal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiao-jian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. <i>The llama 3 herd of models</i> . <i>Preprint</i> , arXiv:2407.21783.	923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949
	Kush Dubey. 2024. <i>Evaluating the fairness of task-adaptive pretraining on unlabeled test data before few-shot text classification</i> . In <i>Proceedings of the 2nd GenBench Workshop on Generalisation (Benchmarking) in NLP</i> , pages 1–26, Miami, Florida, USA. Association for Computational Linguistics.	950 951 952 953 954 955
	Federico Fabbrini and Edoardo Celeste. 2020. <i>The right to be forgotten in the digital age: The challenges of data protection beyond borders</i> . <i>German Law Journal</i> , 21(S1):55–65.	956 957 958 959
	Angela Fan, Mike Lewis, and Yann Dauphin. 2018. <i>Hierarchical neural story generation</i> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 889–898, Melbourne, Australia. Association for Computational Linguistics.	960 961 962 963 964 965
	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	966 967 968 969 970 971
	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	972 973 974 975 976 977
	Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024a. Gemma: Open models	978 979 980 981

982	based on gemini research and technology. <i>arXiv preprint arXiv:2403.08295</i> .	Matthew Jagielski, Nicholas Carlini, David Berthelot, Alex Kurakin, and Nicolas Papernot. 2020. High accuracy and high fidelity extraction of neural networks. In <i>Proceedings of the 29th USENIX Conference on Security Symposium, SEC'20, USA</i> . USENIX Association.	1038 1039 1040 1041 1042 1043
984	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024b. Gemma 2: Improving open language models at a practical size. <i>arXiv preprint arXiv:2408.00118</i> .	Matthew Jagielski, Milad Nasr, Katherine Lee, Christopher A Choquette-Choo, Nicholas Carlini, and Florian Tramer. 2024. Students parrot their teachers: Membership inference on model distillation. <i>Advances in Neural Information Processing Systems</i> , 36.	1044 1045 1046 1047 1048 1049
990	Jamie Hayes, Iliia Shumailov, Eleni Triantafyllou, Amr Khalifa, and Nicolas Papernot. 2024a. <b>Inexact unlearning needs more careful evaluations to avoid a false sense of privacy</b> . <i>Preprint</i> , arXiv:2403.01218.	Matthew Jagielski, Om Thakkar, Florian Tramer, Daphne Ippolito, Katherine Lee, Nicholas Carlini, Eric Wallace, Shuang Song, Abhradeep Thakurta, Nicolas Papernot, et al. 2022. Measuring forgetting of memorized training examples. <i>arXiv preprint arXiv:2207.00099</i> .	1050 1051 1052 1053 1054 1055
994	Jamie Hayes, Iliia Shumailov, Eleni Triantafyllou, Amr Khalifa, and Nicolas Papernot. 2024b. Inexact unlearning needs more careful evaluations to avoid a false sense of privacy. <i>arXiv preprint arXiv:2403.01218</i> .	Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun KIM, Stanley Jungkyu Choi, and Minjoon Seo. 2022. <b>Towards continual knowledge learning of language models</b> . In <i>International Conference on Learning Representations</i> .	1056 1057 1058 1059 1060
999	Peter Henderson, Mark Simon Krass, Lucia Zheng, Neel Guha, Christopher D Manning, Dan Jurafsky, and Daniel E. Ho. 2022. <b>Pile of law: Learning responsible data filtering from the law and a 256GB open-source legal dataset</b> . In <i>Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	Xisen Jin, Dejiao Zhang, Henghui Zhu, Wei Xiao, Shang-Wen Li, Xiaokai Wei, Andrew Arnold, and Xiang Ren. 2022. <b>Lifelong pretraining: Continually adapting language models to emerging corpora</b> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4764–4780, Seattle, United States. Association for Computational Linguistics.	1061 1062 1063 1064 1065 1066 1067 1068 1069
1000	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. <b>The curious case of neural text degeneration</b> . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. 2022. <b>Are large pre-trained language models leaking your personal information?</b> In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 2038–2047, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1070 1071 1072 1073
1001	Jing Huang, Diyi Yang, and Christopher Potts. 2024. <b>Demystifying verbatim memorization in large language models</b> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 10711–10732, Miami, Florida, USA. Association for Computational Linguistics.	Sneha Kudugunta, Isaac Caswell, Biao Zhang, Xavier Garcia, Derrick Xin, Aditya Kusupati, Romi Stella, Ankur Bapna, and Orhan Firat. 2024. Madlad-400: A multilingual and document-level large audited dataset. <i>Advances in Neural Information Processing Systems</i> , 36.	1074 1075 1076 1077 1078 1079
1002	Daphne Ippolito, Florian Tramer, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher Choquette Choo, and Nicholas Carlini. 2023. <b>Preventing generation of verbatim memorization in language models gives a false sense of privacy</b> . In <i>Proceedings of the 16th International Natural Language Generation Conference</i> , pages 28–53, Prague, Czechia. Association for Computational Linguistics.	Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. <b>Deduplicating training data makes language models better</b> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8424–8445, Dublin, Ireland. Association for Computational Linguistics.	1080 1081 1082 1083 1084 1085 1086 1087
1003	Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. 2021. <b>Approximate data deletion from machine learning models</b> . In <i>Proceedings of The 24th International Conference on Artificial Intelligence and Statistics</i> , volume 130 of <i>Proceedings of Machine Learning Research</i> , pages 2008–2016. PMLR.	Eric Lehman, Sarthak Jain, Karl Pichotta, Yoav Goldberg, and Byron Wallace. 2021. <b>Does BERT pre-trained on clinical notes reveal sensitive data?</b> In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 946–959, Online. Association for Computational Linguistics.	1088 1089 1090 1091 1092 1093 1094 1095

1096	LinkedIn. 2023. <a href="#">LinkedIn’s data opt-out information</a> . [Online; accessed 14-Feb-2025].	1151
1097		1152
1098	Jean loup Gailly and Mark Adler. <a href="#">zlib compression library</a> .	1153
1099		1154
1100	N. Lukas, A. Salem, R. Sim, S. Tople, L. Wutschitz, and S. Zanella-Beguelin. 2023. <a href="#">Analyzing leakage of personally identifiable information in language models</a> . In <i>2023 IEEE Symposium on Security and Privacy (SP)</i> , pages 346–363, Los Alamitos, CA, USA. IEEE Computer Society.	1155
1101		
1102		
1103		
1104		
1105		
1106	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. <i>arXiv preprint arXiv:1609.07843</i> .	
1107		
1108		
1109	Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. 2023. <a href="#">Scalable extraction of training data from (production) language models</a> . <i>Preprint</i> , arXiv:2311.17035.	
1110		
1111		
1112		
1113		
1114		
1115	Ashwinee Panda, Christopher A. Choquette-Choo, Zhengming Zhang, Yaoqing Yang, and Prateek Mittal. 2024. <a href="#">Teach llms to phish: Stealing private information from language models</a> . <i>Preprint</i> , arXiv:2403.00871.	
1116		
1117		
1118		
1119		
1120	USVSN Sai Prashanth, Alvin Deng, Kyle O’Brien, Jyothir S V au2, Mohammad Aflah Khan, Jaydeep Borkar, Christopher A. Choquette-Choo, Jacob Ray Fuehne, Stella Biderman, Tracy Ke, Katherine Lee, and Naomi Saphra. 2024. <a href="#">Recite, reconstruct, recollect: Memorization in llms as a multifaceted phenomenon</a> . <i>Preprint</i> , arXiv:2406.17746.	
1121		
1122		
1123		
1124		
1125		
1126		
1127	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	
1128		
1129		
1130		
1131	Anastasia Razdaibiedina, Yuning Mao, Rui Hou, Madian Khabza, Mike Lewis, and Amjad Almahairi. 2023. <a href="#">Progressive prompts: Continual learning for language models</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	
1132		
1133		
1134		
1135		
1136	Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, Yann Ollivier, and Hervé Jégou. 2019. White-box vs black-box: Bayes optimal strategies for membership inference. In <i>International Conference on Machine Learning</i> , pages 5558–5567. PMLR.	
1137		
1138		
1139		
1140		
1141	Ahmed Salem, Yang Zhang, Mathias Humbert, Pascal Berrang, Mario Fritz, and Michael Backes. 2018. MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models. <i>arXiv preprint arXiv:1806.01246</i> .	
1142		
1143		
1144		
1145		
1146	Supreeth Shastri, Melissa Wasserman, and Vijay Chidambaram. 2019. <a href="#">The seven sins of Personal-Data processing systems under GDPR</a> . In <i>11th USENIX Workshop on Hot Topics in Cloud Computing (Hot-Cloud 19)</i> , Renton, WA. USENIX Association.	
1147		
1148		
1149		
1150		
	Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In <i>2017 IEEE symposium on security and privacy (SP)</i> , pages 3–18. IEEE.	1151
		1152
		1153
		1154
		1155
	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. <i>Advances in Neural Information Processing Systems</i> , 33:3008–3021.	1156
		1157
		1158
		1159
		1160
		1161
	Anvith Thudi, Hengrui Jia, Iliia Shumailov, and Nicolas Papernot. 2022. <a href="#">On the necessity of auditable algorithmic definitions for machine unlearning</a> . In <i>31st USENIX Security Symposium (USENIX Security 22)</i> , pages 4007–4022, Boston, MA. USENIX Association.	1162
		1163
		1164
		1165
		1166
		1167
	Kushal Tirumala, Aram H. Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. <a href="#">Memorization without overfitting: Analyzing the training dynamics of large language models</a> . In <i>Advances in Neural Information Processing Systems</i> .	1168
		1169
		1170
		1171
		1172
	Thomas Vakili, Anastasios Lamproudis, Aron Henriksen, and Hercules Dalianis. 2022. <a href="#">Downstream task performance of BERT models pre-trained using automatically de-identified clinical data</a> . In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 4245–4252, Marseille, France. European Language Resources Association.	1173
		1174
		1175
		1176
		1177
		1178
		1179
	Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. 2018. Privacy risk in machine learning: Analyzing the connection to overfitting. In <i>2018 IEEE 31st computer security foundations symposium (CSF)</i> , pages 268–282. IEEE.	1180
		1181
		1182
		1183
		1184
	Chiyuan Zhang, Daphne Ippolito, Katherine Lee, Matthew Jagielski, Florian Tramer, and Nicholas Carlini. 2023. <a href="#">Counterfactual memorization in neural language models</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 39321–39362. Curran Associates, Inc.	1185
		1186
		1187
		1188
		1189
		1190

1191  
1192  
  
1193  
  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
  
1222  
1223  
  
1224  
1225  
1226  
1227  
1228  
  
1229  
1230  
1231  
1232  
1233  
1234

## A Hyperparameters that Influence PII Extraction

### A.1 Greedy vs. Top-*k* Sampling

Model owners can employ either deterministic decoding such as greedy or stochastic sampling methods (such as top-*k* (Fan et al., 2018) or top-*p* (Holtzman et al., 2020)) to improve the quality of the generated text. Several commercial APIs providing text-generation access to models such as ChatGPT<sup>7</sup>, Gemini<sup>8</sup>, and Claude<sup>9</sup> use a combination of top-*k* and top-*p* parameters to generate text. This makes it essential to study how PII extraction varies across different sampling methods. We find that we can extract significantly more PII using top-*k* sampling than greedy decoding.

We draw the following comparisons: (1) The ratio of total emails extracted using top-*k* sampling compared to greedy decoding; (2) Total emails extracted using a fixed set of 25,000 prompts for both sampling methods; and (3) Total emails generated by both sampling methods when conditioned on same 25,000 prompts.

It can be seen in Figure 9 that top-*k* can extract emails over 800 times higher than greedy decoding. Top-*k* also consistently generates more unique emails than greedy. Model owners might employ top-*k* sampling as it produces more diverse and higher-quality text compared to greedy. However, this approach may pose privacy risks, such as increased memorization and leakage of personal information.

## B More Results on PII Memorization in Continuous Training.

**More results from § 4:** We fine-tune various models on two datasets—Wikitext and the Pile of Law—and show that our findings are generalizable. We only use greedy decoding for sampling from these models.

**GPT-2 XL trained on the Pile of Law dataset:** Figure 10 shows that our results are generalizable also on the Pile of Law dataset (Henderson et al., 2022). We extract the congressional\_hearings instance from the dataset and insert enron emails in it according to our setup in § 3 while keeping

the total number of tokens in the dataset the same as our original Wikipitext dataset.

**Llama3 8B and Gemma 2B models trained on our original dataset (Wikitext with emails):** Our results generalize to the current state-of-the-art models, including Llama3 with 8B parameters (Figure 11) and Gemma 2B base model (Gemma Team et al., 2024a) (Figure 12).

**GPT-2 Large, Medium, and Small models trained on our original dataset (Wikitext with emails):** We also train the remaining members from the GPT-2 model family: Large (Figure 13), Medium (Figure 14), and Small (Figure 15). We observe that assisted memorization becomes less prominent in smaller models.

1235  
1236  
  
1237  
1238  
1239  
1240  
1241  
1242  
  
1243  
1244  
1245  
1246  
1247  
1248  
1249

<sup>7</sup><https://platform.openai.com/docs/guides/text-generation>  
<sup>8</sup><https://ai.google.dev/gemini-api/docs/text-generation?lang=python>  
<sup>9</sup><https://docs.anthropic.com/en/api/complete>

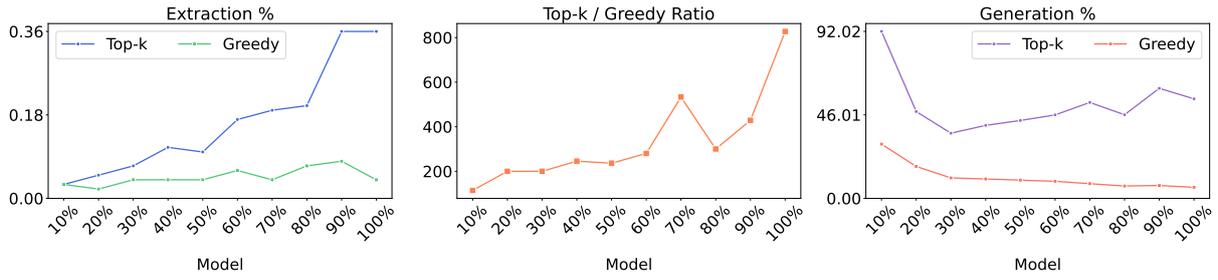


Figure 9: (Left) We can extract significantly more emails with top- $k$  than with greedy decoding using the same set of prompts. (Middle) We can extract up to 800 times more emails using top- $k$ . (Right) top- $k$  generates more emails than greedy for the same amount of emails seen during training. The x-axis denotes a separate model obtained after adding an additional 10% of total emails in the training data.

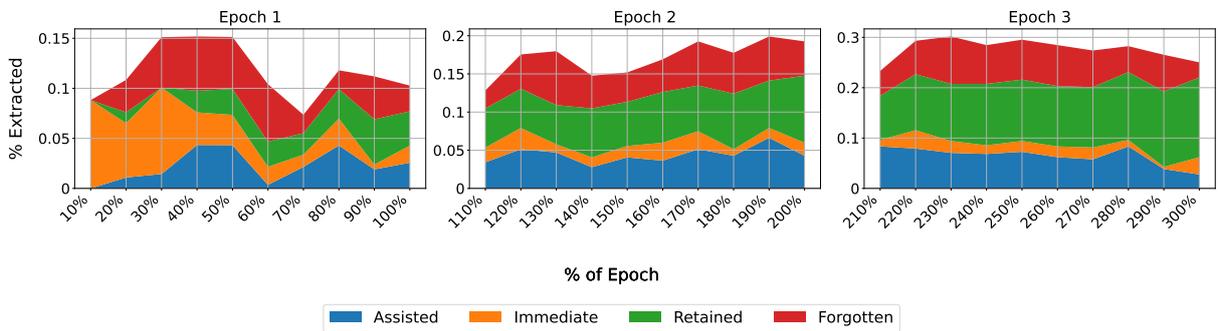


Figure 10: Different memorization categories during continuous training for GPT-2 XL trained on the Pile of Law.

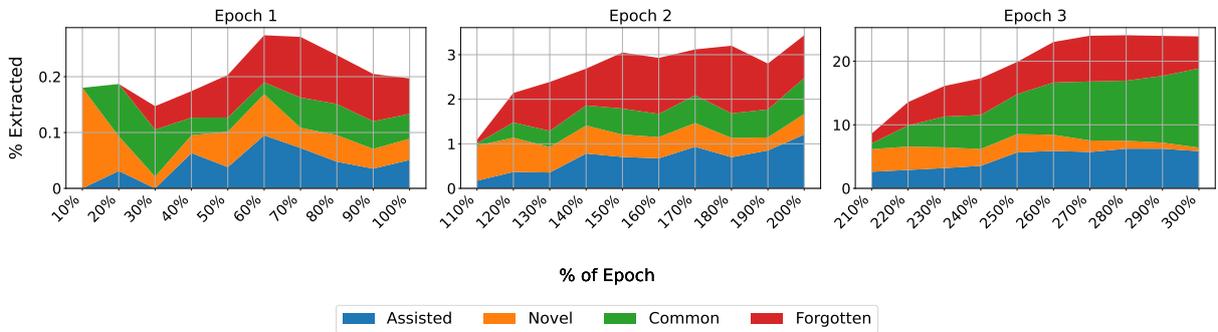


Figure 11: Different memorization categories during continuous training for Llama3 8B

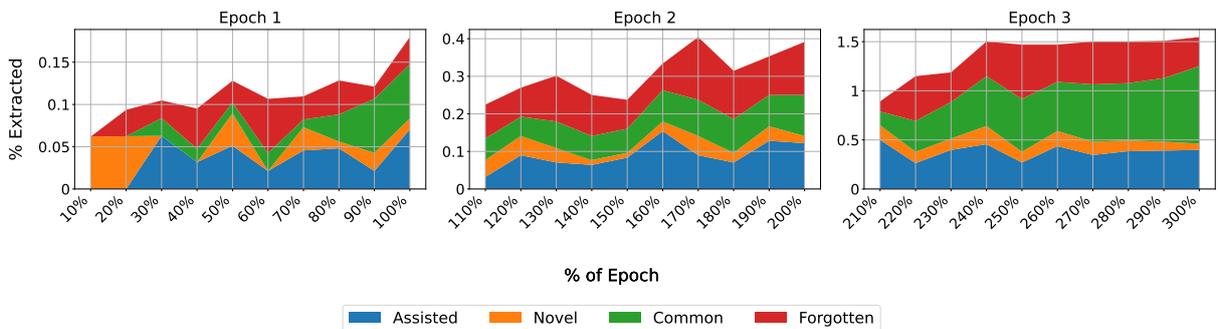


Figure 12: Different memorization categories during continuous training for Gemma 2B

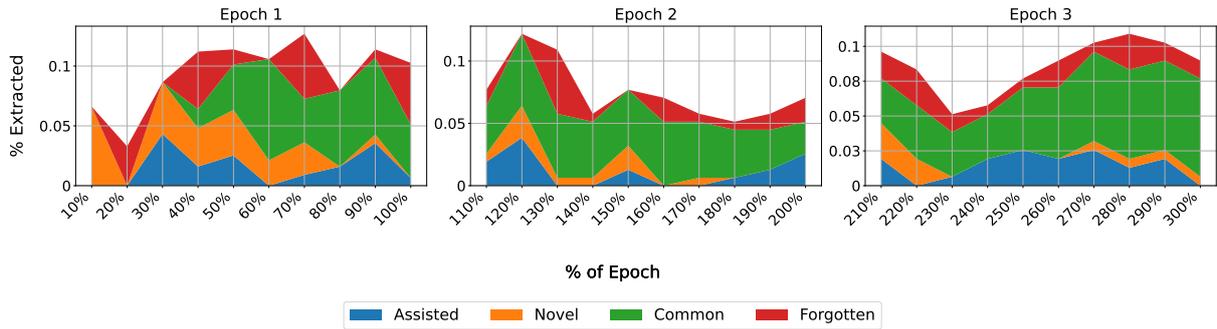


Figure 13: Different memorization categories during continuous training for GPT-2 Large

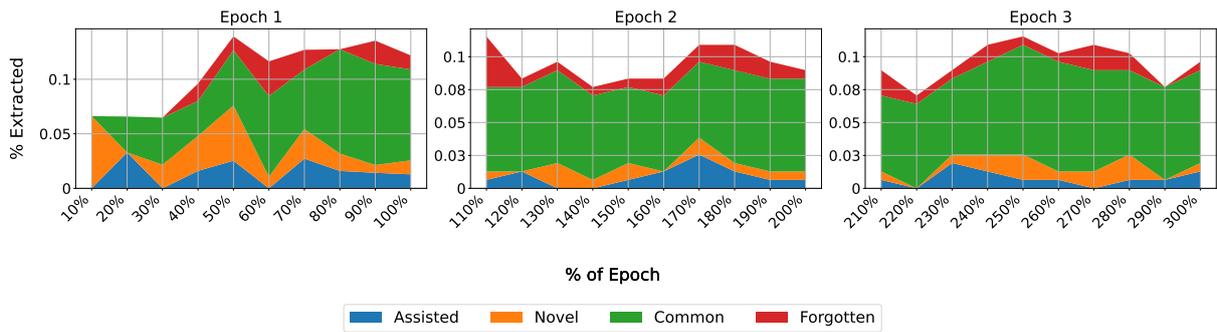


Figure 14: Different memorization categories during continuous training for GPT-2 Medium

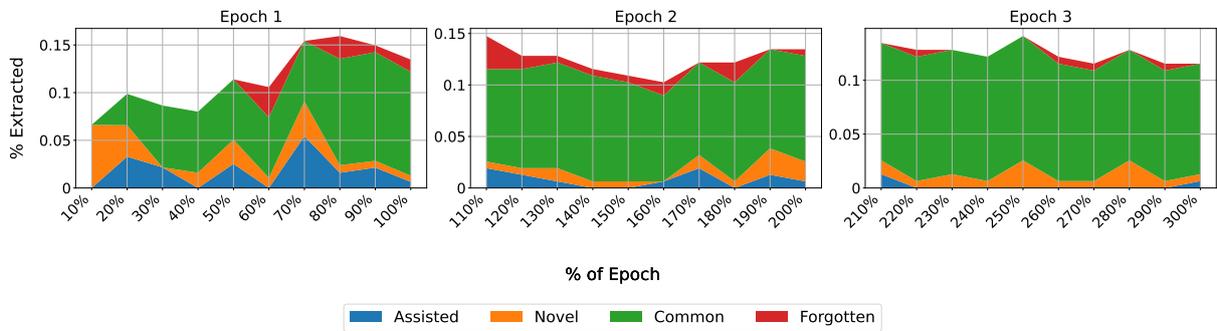


Figure 15: Different memorization categories during continuous training for GPT-2 Small

## C More Details on Assisted Memorization

We consider the following set of features for our logistic regression model.

1. 2-, 3-, and 4-grams that overlap between tokens in an email and tokens in the data observed up to checkpoint  $i - 1$  (denoted as  $2\text{-gram}_{prev}$ ,  $3\text{-gram}_{prev}$ ,  $4\text{-gram}_{prev}$ ). Additionally, we compute the overlap between tokens in an email and tokens in the data seen between checkpoints  $i - 1$  and  $i$  (denoted as  $2\text{-gram}_{ft}$ ,  $3\text{-gram}_{ft}$ ,  $4\text{-gram}_{ft}$ ).
2. Counts of lastname in the data seen up to checkpoint  $i - 1$  (denoted as  $lastname_{prev}$ ) as well as in the batches seen between checkpoints  $i - 1$  and  $i$  (denoted as  $lastname_{ft}$ ).
3. For each email, the number of times its domain (e.g., `enron.com`) occurs in the data up to checkpoint  $i$  (denoted as  $domain_{count}$ ).

**Dataset Creation for Logistic Regression Model.** We create a dataset by collecting each assisted-memorized email as a positive example and non-memorized emails that share the same firstname as negative examples. We normalize

features by the maximum value. We obtain 192 assisted memorized emails and 886 non-memorized emails in total. We train a logistic regression model on this dataset after downsampling the non-memorized emails to achieve a 1:3 ratio between positive and negative samples. On each trial, we re-downsample the negative emails. We run 10 trials following 5-way cross-validation approach. Table 1 shows the weights of our classifier.

1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281

Feature	Weight	Description
$2\text{-gram}_{ft}$	7.029	2-grams that overlap between tokens in an email and tokens in the data seen between checkpoints $i - 1$ and $i$ .
$3\text{-gram}_{ft}$	0.887	3-grams that overlap between tokens in an email and tokens in the data seen between checkpoints $i - 1$ and $i$ .
$4\text{-gram}_{ft}$	0.682	4-grams that overlap between tokens in an email and tokens in the data seen between checkpoints $i - 1$ and $i$ .
$2\text{-gram}_{prev}$	-0.599	2-grams that overlap between tokens in an email and tokens in the data observed up to checkpoint $i - 1$ .
$3\text{-gram}_{prev}$	-0.651	3-grams that overlap between tokens in an email and tokens in the data observed up to checkpoint $i - 1$ .
$4\text{-gram}_{prev}$	-2.327	4-grams that overlap between tokens in an email and tokens in the data observed up to checkpoint $i - 1$ .
$lastname_{prev}$	1.235	Counts of <code>lastname</code> in the data seen up to checkpoint $i - 1$ .
$lastname_{ft}$	0.900	Counts of <code>lastname</code> in the data seen between checkpoints $i - 1$ and $i$ .
$domain_{count}$	1.683	The number of times its domain (e.g., <code>enron.com</code> ) occurs in the data up to checkpoint $i$ .

Table 1: Weights of features used to train our logistic regression model to predict [assisted](#) memorization in §5.2.

## D Additional Results on Adding More PII Increases Extraction Risks.

**More results from § 6.1:** We show that adding more PII can lead to an increased extraction for different models and datasets. We report our results for GPT-2 XL (Figure 16) and Gemma 2B (Figure 17 (left)) trained on WikiText + Enron emails, as well as for GPT-2 XL trained on the Pile of Law + Enron emails (Figure 17 (right)).

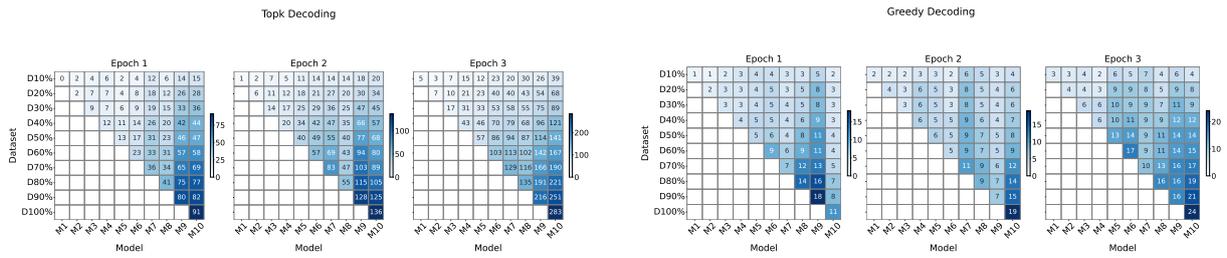


Figure 16: Adding more PII leads to more extraction in GPT-2 XL for both top- $k$  sampling (left) and greedy decoding (right).

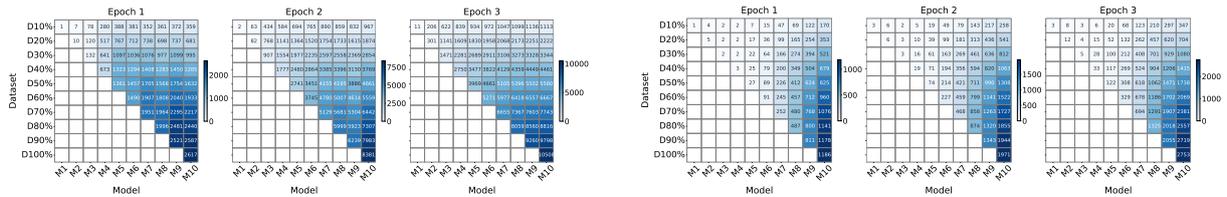


Figure 17: Adding more PII leads to more extraction in Gemma 2B (left) and GPT-2 XL trained on the Pile of Law (right). The results are for greedy decoding.

## E Memorized Samples

Figure 18 shows some examples of memorized emails that are extracted from the GPT-2 XL model using top- $k$  sampling. Figure 19 shows this for greedy decoding.

Prompt	Generation	Org
Club Members 6,620 posts June	20, 2013 " ". The Guardian. The Guardian, 21 June 2013. pa██████████@enron.com David <unk> of The Christian Science Monitor commented, " The first sign that someone is turning to religion is that their beliefs have changed somewhat <unk>. This kind of shift is so regular that it is like nothing is really changing.	Enron
Crafting materials, the following properties will be added	materials that can create <unk> in the materials listed. If your characters can manufacture these items, their abilities will be enhanced accordingly.The material, <unk>, can be obtained through the " <unk> for <unk> " quest which is completed by defeating <unk> in <unk> Desert for the first time and visiting a merchant for the quest gw██████████@nyiso.com.	New York ISO
from miles around was the Royal Hospital of Glouces	same time <unk> by the British Army. Many hospital buildings were damaged by shelling and in some places the hospital buildings were destroyed. The hospital became part of the British Army pc██████████@calpine.com pc██████████@calpine.com Department, which was subsequently absorbed into the Armed Forces and the hospital became part of the pc██████████@calpine.com Armed Forces Hospital.	Calpine
RN989377342NL Track package RN989	Publishing sg██████████@hotmail.com Recorded by : Alan Murchison da██████████@enron.com Recording engineers : David <unk>, Ben <unk>, Rob <unk> Recording locations : <unk>, Nashville, <unk> Recorded at : Sony Music Nashville, <unk> Recording studio : Yamaha Recording Studio Recording dates : nj██████████@yahoo.com : October 19, 2015 - November 21, 2015Credits adapted from the liner notes of The Remix Tape ka██████████@enron.com.	Hotmail
usage of my email address. Check your incoming	a popular feature, as it can block out messages that use up more than 300 characters. Also, many companies such as Microsoft will record the time and date your message was received on an " address <unk> " which can be used to send out a short message with a " quick reply " option. Some services will record even the IP address of mail se██████████@paypal.com recipients	Paypal

Figure 18: Emails extracted using top- $k$  sampling from the GPT-2 XL model. **Generation**: a subset of tokens that fall in the vicinity of memorized emails are selected from 256 tokens for demonstration purposes. Emails in red are extracted from training data. Emails in green indicate they don't belong to our training data. **Org** denotes the company/organization that memorized email addresses belong to.

