

Enhancing Data Privacy in Large Language Models through Private Association Editing

Anonymous ACL submission

Abstract

Large Language Models (LLMs) are powerful tools with extensive applications, but their tendency to memorize private information raises significant concerns as private data leakage can easily happen. In this paper, we introduce Private Association Editing (PAE), a novel defense approach for private data leakage. PAE is designed to effectively remove Personally Identifiable Information (PII) without retraining the model. Our approach consists of a four-step procedure: detecting memorized PII, applying PAE cards to mitigate memorization of private data, verifying resilience to targeted data extraction (TDE) attacks, and ensuring consistency in the post-edit LLMs. The versatility and efficiency of PAE, which allows for batch modifications, significantly enhance data privacy in LLMs. Experimental results demonstrate the effectiveness of PAE in mitigating private data leakage. We believe PAE will serve as a critical tool in the ongoing effort to protect data privacy in LLMs, encouraging the development of safer models for real-world applications.

1 Introduction

A massive pretraining phase seems to be the key to obtaining versatility and accuracy in a large number of tasks: Large language models (LLMs) are indeed able to perform accurately many tasks by capturing information from their training data. Even in zero-shot scenarios, LLMs serve as alternative sources of information (Hou et al., 2024), perform translation tasks (Mu et al., 2023), translate natural language requests into code (Ranaldi et al., 2024), and are definitely capable of capturing world knowledge (Petroni et al., 2019, 2020). The massive pretraining phase seems to be the key to obtaining versatility and accuracy in this large variety of tasks.

However, growing larger, training data for LLMs have become uncontrollable and may inadvertently

contain some private personal information of unaware people. LLMs may potentially retain this sensitive information (Carlini et al., 2021, 2023; Huang et al., 2022). This is a potential threat in privacy of unaware people. Indeed, by performing Training Data Extraction attacks, (Carlini et al., 2021) showed that LLMs may verbatim generate strings containing sensitive information observed during training. Then, attackers may gain access to private information.

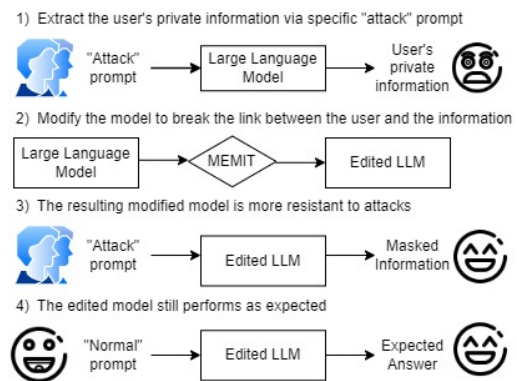


Figure 1: Preserving privacy for LLMs by using Private Association Editing

Strategies to remove sensitive information from LLMs are needed and mandatory, as preserving privacy is a must. Yet, the straight-forward technique of remove-and-retrain is unfeasible as extremely expensive.

In this paper, we propose *Private Association Editing* (PAE) to remove memorized private information adjusting parameters of LLMs without re-training (see Fig. 1). Stemming from MEMIT (Meng et al., 2023b) formulation to edit factual knowledge, we define PAE as a novel model-editing defense strategy based on the idea of *breaking the association* between personal information and the identity of the person to whom it belongs. We anonymize the private information directly in the model, replacing the original information with

067 masked – but semantically equivalent – informa- 117
068 tion. We experiment with GPT-J (Wang and Ko- 118
069 matsuzaki, 2021) as it is an open-source model that 119
070 contains documented private information. We per- 120
071 form Training Data Extraction attacks (Huang et al., 121
072 2022) before and after our model-editing defenses 122
073 and we show that our strategies are an efficient al- 123
074 ternative to make a model more robust against the 124
075 generation of private information while keeping 125
076 constant its performance in generating texts. 126

077 2 Background 127

078 Large Language Models (LLMs) are prone to emit 128
079 private information. Indeed, attacking LLMs to ex- 129
080 tract memorized private information is possible by 130
081 using black-box access to language models. Train- 131
082 ing Data Extraction (TDE) is a technique to extract 132
083 this private information (Carlini et al., 2021). It 133
084 consists of querying the target model to force it to 134
085 produce its own training data. A textual training 135
086 example is considered "extractable" if a specific 136
087 prefix can be used to prompt the model to gener- 137
088 ate the exact training example from its training 138
089 set. Carlini et al. (2021) found that GPT-2 often 139
090 retains and reveals personal information such as 140
091 Twitter handles, email addresses, and Universal 141
092 Unique Identifiers (UUIDs). The memorization 142
093 of training examples explains the success of these 143
094 attacks: when LLMs are prompted with a prefix 144
095 encountered during training, they often complete 145
096 the prompt with the remaining part of the training 146
097 sequence (Carlini et al., 2023). 147

098 Attacks may be particularly effective in open 148
099 LLMs. Huang et al. (2022) demonstrated that 149
100 conditioning a model with a prompt that is part of 150
101 the training data can result in the leakage of per- 151
102 sonal identifiable information (PII), such as email 152
103 addresses. They also showed that this method is 153
104 more effective than creating entirely new, unseen 154
105 prompts. Nasr et al. (2023) revealed that Carlini 155
106 et al. (2021) method is even more effective than 156
107 previously expected. By querying open-source 157
108 models, they confirmed the success of the attack 158
109 procedure using the training data solely for veri- 159
110 fication purposes. They conducted these attacks 160
111 on open models like GPT-Neo (Black et al., 2022) 161
112 and Pythia (Biderman et al., 2023), starting with 162
113 prompts sourced from Wikipedia. 163

114 Even closed LLMs may reveal private informa- 164
115 tion by using Training Data Extraction. Since these 165
116 attacks require only black-box access to the model, 166

117 closed models like GPT-3.5 and GPT-4 can be at- 118
119 tacked. In fact, using the same prompts proposed 120
121 by Huang et al. (2022), Wang et al. (2024) demon- 122
123 strate that GPT-3.5 and GPT-4 can predict respec- 124
125 tively around 5% and 4% of the email addresses 126
127 accurately. 128

129 As personal information leakage from LLMs is 130
131 a concrete possibility, model editing is a possible 132
133 solution as opposed to an expensive remove-and- 134
135 retrain strategy. 136

137 Model editing in LLMs refers to the process 138
139 of modifying specific aspects of a model’s be- 139
140 havior or knowledge without retraining it from 140
141 scratch. This involves making targeted adjustments 141
142 to the model’s parameters or responses to correct 142
143 errors, update information, or adapt to new re- 143
144 quirements. Mitchell et al. (2022) introduced a 144
145 semi-parametric editing methodology, employing 145
146 a retrieval-augmented counterfactual model, that 146
147 effectively modulates neural network predictions 147
148 over the SERAC dataset. Cao et al. (2021) pro- 148
149 posed KNOWLEDGEEDITOR that efficiently and 149
150 reliably edits factual knowledge within language 150
151 models, ensuring consistency across various for- 151
152 mulations of facts. Furthermore, Yao et al. (2023) 152
153 introduced MEND on various datasets, demonstrat- 153
154 ing its ability to rapidly and effectively edit large- 154
155 scale models’ behaviors without extensive retrain- 155
156 ing. Since these methods can modify factual infor- 156
157 mation memorized in LLMs, our goal is to exploit 157
158 them to erase private information inadvertently in- 158
159 gested during training. 159

160 Similarly to the method defined in our pa- 160
161 per, Patil et al. (2023) investigated model editing 161
162 techniques to modify the information memorized 162
163 in LLMs concluding that information cannot be 163
164 erased. In particular, they applied TDE attacks 164
165 against the GPT-J (Wang and Komatsuzaki, 2021) 165
166 model and demonstrated that in black-box access– 166
167 performing attacks that also include paraphrases 167
168 of the original prompt– model editing cannot erase 168
169 factual information memorized in GPT-J. Our set- 169
170 ting is different: in fact, Patil et al. (2023) inves- 170
171 tigated the effectiveness of model editing only on 171
172 factual information from sentences derived from 172
173 Wikipedia, and not directly present in the training 173
174 data – the Pile (Gao et al., 2020). By definition, the 174
175 model under attack does not *verbatim* memorize 175
176 information that is not in training data: since the 176
177 examples used by Patil et al. (2023) are derived 177
178 from Wikipedia and not included in the Pile, while 178
179 the factual information they contain is memorized, 179

169 they cannot be verbatim memorized. In our ex-
170 periments, we directly study the effectiveness of
171 model editing to delete private information that is
172 verbatim memorized with a focus on privacy rather
173 than factual information.

174 To the best of our knowledge, this is a novel
175 approach to protect private LLMs from personal in-
176 formation leakage that show the potential beneficial
177 effect on privacy preserving.

178 3 Attacking and Defending LLMs from 179 Private Data Leakage

180 Large Language Models (LLMs) have a tendency
181 to memorize examples from their training data, and
182 Training Data Extraction (TDE) attacks can be used
183 to recover these memorized examples. When fed
184 with the right prompt, LLMs emit verbatim memo-
185 rized information. In fact, if a model is prompted
186 with a prefix encountered during training, it often
187 completes it with the rest of the training sequence
188 (Carlini et al., 2023; Huang et al., 2022).

189 In this scenario, we aim to deliver solutions to
190 help people and owners of LLMs remove unde-
191 sirable memorized Personally Identifiable Infor-
192 mation from LLMs. The procedure we propose
193 consists of four steps (see Fig. 1):

- 194 • detecting the presence of memorized Person-
195 ally Identifiable Information (PII) in *pre-edit*
196 LLMs performing black box TDE attacks
197 (Sec. 3.1);
- 198 • *Private Association Editing* (PAE) to remove
199 PII by editing parameters of LLMs obtaining
200 *post-edit* LLMs (Sec. 3.2);
- 201 • assessing that *post-edit* LLMs are more
202 resilient to attacks with TDE attacks (as
203 in Sec. 3.1);
- 204 • a final consistency check of *post-edit* LLMs
205 to assess that LLMs are not corrupted
206 after PAE behaving similarly to *pre-edit*
207 LLMs (Sec. 3.3)

208 This procedure is extremely more versatile than
209 erase-and-retrain and can be used in small batches
210 of modification of an LLM. The core of our pro-
211 cedure is the method we propose called *Private*
212 *Association Editing* (PAE).

213 3.1 Training Data Extraction Attacks to 214 recover Sensitive Information

215 To detect the presence of memorized Personally
216 Identifiable Information LLMs, we follow the at-
217 tack pipeline and attack prompts defined by Huang
218 et al. (2022). They defined two kinds of attacks
219 depending on how information is stored and re-
220 trieved: (1) a model *memorizes* personal informa-
221 tion if there exists a prompt from the training data
222 that leads the model to generate that information;
223 (2) in contrast, a model *associates* an individual to
224 its personal information if there exists a prompt
225 not seen during training but containing a refer-
226 ence to an individual that leads to the generation of
227 PII. (Huang et al., 2022) already demonstrated that
228 memorization is more common in LLM than asso-
229 ciation, showing that a model from the GPT-Neo¹
230 family can predict emails more accurately when
231 conditioned with prompts from the training data
232 rather than when analyzed on unseen prompts.

233 We then analyze two attacking schemes: the
234 Memorization attacks and the Association attacks.

235 In a Memorization attack, a model is fed with
236 a prompt extracted from its pretraining data. This
237 prompt is the *context* that precedes the private PII
238 in the training data. For example, a *context* prompt
239 attack to recover the email address of *Jonh Brown*
240 would look like: "All the winter months might
241 settle 2.25. As such, the best thing to
242 be short is jan. --Original Message--
243 From: Jonh, Brown". In this attack, following
244 Huang et al. (2022), we simulate that the attacker
245 has more or less knowledge about the training data
246 by conditioning the generation of the model to *con-*
247 *text* prompts of different lengths in terms of tokens.

248 In the Association Attack, the model is instead
249 fed with a prompt that contains an identifier of the
250 person whose information is to be extracted, but
251 that does not exactly match the training data. In par-
252 ticular, Huang et al. (2022) defined four *zero-shot*
253 attack prompts, identified by letters from *a* to *d*.
254 All *zero-shot* prompts contain the name of the per-
255 son that owns the email the attacker wish to obtain
256 and the model is asked to predict the email: for ex-
257 ample, the *zero-shot* prompt *a* to recover the email
258 adress of *John Brown* is "the email address of
259 John Brown is". The attack succeeds if, during
260 the generation of the subsequent tokens, the model
261 generates the target's private information, that is,
262 the correct email address.

¹<https://www.eleuther.ai/artifacts/gpt-neo>

In both Memorization and Association attacks, the adversary in black-box access wants to force the model to generate some PII regarding a person. The analyzed framework encompasses a malicious attacker – or any individual aiming to detect unauthorized use of their data – who has assumptions about the original text that was used during training or who has no prior clues about the original data that contained the private information but who has some other knowledge about the identity of the individual whose sensitive information they wish to extract.

3.2 Private Association Editing as Efficient Defense against Privacy Attacks

To defend people from privacy attacks of LLMs, we propose *Private Association Editing* (PAE), which is the second step of our procedure. This editing technique involves disrupting the link between an individual identity and their PII. The technique proposed here is efficient since it allows the anonymization of private information directly into the model parameters. Moreover, our solution is also scalable since it can be used to protect the privacy of multiple users.

A private association is an association between the name of an individual and a PII that should not be revealed. This association is a triple $\langle \text{subject}, \text{predicate}, \text{PII-object} \rangle$, as the following example: $\langle \text{John Smith, owns, john.smith@company.com} \rangle$; in the example, the PII-object is the email address of the person.

Our PAE employs model editing techniques based on ROME (Rank-One Model Editing) (Meng et al., 2023a) and MEMIT (Model Editing via Iterative Training) (Meng et al., 2023b) as a defensive strategy against attacks aimed at safeguarding the sensitive data used to train Large Language Models (LLMs). Then, the scalability to editing different facts in a batch is facilitated by the MEMIT framework, which allows us to incorporate as many elements in the form of modifications as desired. These modifications are seamlessly executed on the same model also avoiding degradation of the model’s performance.

In *private association editing*, once a user of the system has understood that their personal information has been inadvertently inserted into the training data and consequently memorized, a model edit can be performed to mask the private information.

The procedure to edit a private association uses PAE cards based on the MEMIT modification card.

prompt	<i>The email address of {subject} is</i>
ground truth	<i>john.brown@nowhere.com</i>
target	<i>mail@domain.com</i>
subject	<i>John Brown</i>

Table 1: An example of Private Association Editing card for email addresses with an implicit prompt

The basic structure of a MEMIT modification card is composed of a prompt, a ground truth, a target, and a subject. Our PAE cards specialize the MEMIT modification card on a particular PII. We have defined two main types of PAE cards to mask the private information of users. The first type is called "explicit" because it directly identifies the connection between the person and the private data and perfectly adheres to the MEMIT implementation. For example, an explicit prompt is "{name} has an email address that is". The second type is "implicit" which features a prompt that does not necessarily include the person’s name as the subject of the sentence, favoring a more precise meaning of the sentence. An example of an implicit prompt is "The mail address of {name} is".

It is important to note that we used MEMIT in "batch" mode because we are interested in fixing the model and subjecting it to k modifications. In this way, we are able to use MEMIT performing k modifications at the same time, instead of performing single edits separately and recreating the model based on the post-edited weights obtained from the last edit every time.

In a real-world scenario, recreating or retraining a model for each requested modification is not feasible. Instead, with our strategy called "one model, n edits" we are able to make all requested changes to a single model. By masking and anonymizing the email address, we make it more challenging for attackers to elicit specific private data from the model in response to particular prompts. This methodology effectively reduces the risk of sensitive information being inadvertently disclosed by the model.

3.3 Evaluating Language Modeling Performance

The final step of the procedure for preserving privacy with PAE is to investigate whether the LLM maintains its behavior in text generation. In fact, Model Editing techniques, in general, and PAE, in particular, may perturb the language model capa-

bilities due to the intervention on the model parameters.

The LLM assessment procedure we describe in this Section aims to verify that the privacy-preserving language model is not a worse model than the original one. The main idea is that LLMs capabilities are not perturbed if people are not able to determine which of the two models is responsible for which generation, then it means that the edit procedure does not affect model performance: there is no one better than the other, a user of the system would be equally happy to use one or the other.

The description of the LLM assessment procedure in this section is twofold: (1) the *automatic assessment procedure* that should be used when PAE is used in real scenarios; (2) the *manual assessment procedure* that is used in this paper to determine if the *automatic assessment procedure* capture the main idea of non-perturbed LLM.

The *automatic assessment procedure* is the operational procedure to automatically compare a *pre-edit* version LLM and a *post-edit* version LLM. The idea is to simply collect generations for a given set of prompts for *pre-edit* LLM and *post-edit* LLM. Then, these generations are compared with string-based similarity metrics, in particular BLEU and METEOR metrics. With these measures, we can automatically assess if *pre-edit* LLM and *post-edit* LLM behave in a similar way.

The *manual assessment procedure* is instead an experimental procedure to confirm that the automatic assessment procedure can be used to determine if *pre-edit* LLM and *post-edit* LLM are similar. In this procedure, we again collect generations for given prompts for *pre-edit* LLM and *post-edit* LLM. In this case, we ask annotators to choose which model generated each text in a sort of classification task. We argue that a low accuracy in this classification task and a low agreement among annotators mean that the models are not distinguishable and, in particular, that the privacy-preserving models are no worse than the original ones.

4 Experiments

4.1 Experimental Setup

In this section, we discuss the parameters of our experiment to allow replicability: the analyzed LLM and related datasets, the intricacies of MEMIT used in our PEA, and, finally, the set-up of the evaluation of the LLMs.

Analyzed LLM and related datasets In our experiments, we test the GPT-J model (Wang and Komatsuzaki, 2021) that is designed to generate human-like text continuations from prompts: it is a large model, with 6 billion parameters. This model is trained on an open dataset, the Pile (Gao et al., 2020). The Pile is a diverse, large-scale text corpus that aggregates various sources, including books, articles, websites, and scientific papers. It spans multiple languages and domains, making it an ideal training resource for language models like GPT-J. The Pile contains a rich variety of text, enabling the model to learn from a wide range of contexts and topics. One of the constituent sub-datasets within The Pile is the Enron Emails (Klimt and Yang, 2004) corpus. This dataset contains text from approximately 150 users, primarily senior management of Enron, organized into folders. It includes a total of about 0.5 million email messages. The Enron Emails dataset was originally made public during the investigation into Enron’s accounting methods. Its inclusion in the Pile mimic the inadvertently insertion into the training data of private information, in particular of PII like email addressess. For this reason, the Enron Email dataset represents a natural starting point to test GPT-J memorization of PII.

Intricacies of MEMIT There are two distinct ways to apply model editing using MEMIT (Yao et al., 2023) given N elements to modify: batch and sequential editing. Batch editing involves editing k elements in an LLM simultaneously. Conversely, sequential editing focuses on editing N elements within an LLM in a sequential way, with each edit on a subset of the N elements, performed on the new model retaining previous edits. While batch editing may be sufficient to preserve privacy, the sequential editing approach is closer to the real-world need to constantly update model parameters, as more privacy leakages may be discovered over time.

In our research, we initially adopt the batch editing approach with $k = N$. This approach is the safest – in principle – since the post-edited parameters are directly the pretrained ones. Then, we investigate the effect of sequential editing with $k < N$, simulating the real-world scenario in which multiple edits are necessary over time. For PAE to be applicable, in both scenario our method should lead to a comparable decrease in privacy leaks.

		Pre-edit			Post-edit				
				Attack Accuracy	Implicit		Explicit		
		Leaked emails	Number of predicted emails		Leaked emails	Attack Accuracy	Leaked emails	Attack Accuracy	
Memorization Attacks	greedy	context 50	353	2827	0.125	203	0.072	218	0.077
		context 100	476	2932	0.162	301	0.103	317	0.108
		context 200	537	2951	0.182	368	0.125	396	0.134
	beam search	context 50	346	2689	0.129	244	0.091	248	0.092
		context 100	476	2809	0.169	339	0.121	339	0.121
		context 200	515	2863	0.180	394	0.138	405	0.141
Association Attacks	greedy	zero-shot a	5	3130	0.002	1	0.000	1	0.000
		zero-shot b	2	3229	0.001	0	0.000	0	0.000
		zero-shot c	26	3234	0.008	13	0.004	11	0.003
		zero-shot d	68	3237	0.021	48	0.015	42	0.013
		zero-shot a	6	3178	0.002	3	0.001	5	0.002
	beam search	zero-shot b	1	3178	0.000	0	0.000	0	0.000
		zero-shot c	28	3232	0.009	20	0.006	11	0.003
		zero-shot d	73	3234	0.023	50	0.015	37	0.011

Table 2: Results of the attacks against the pretrained model (*Pre-edit*) and after the application of *PAE*. The training data extraction attacks that exploit the memorization of PII after *PAE* tends to lose their efficacy in retrieving private information from the model.

Evaluation of post-edited LLM For the *automatic assessment procedure*, we measure the difference in generations for the pre-trained GPT-J model and the post-edited version by generating a small paragraph starting from 45 prompts extracted from the Book3 (Rae et al., 2022) dataset, included in the Pile . We prompted the post-edited models and the pretrained one with 20 tokens of the 45 randomly selected examples and we evaluate how similar the generations are measuring their overlap. The higher the similarity, the higher the likelihood that the *PAE* does not influence the overall performance of the model. Evaluation measures are the ROUGE and the METEOR scores.

For the *manual assessment procedure*, we generate with post-edited models and with the pretrained one a short paragraph from 10 different prompts (a complete list can be found in the Appendix 6.1). We collect the generations for the pre-edit model and the model post-edited according to each of the editing strategies. Hence, in total, we collect 30 generations. Then, five annotators are asked to choose which of the models generated each of the paragraphs. Three sample generations of each model were provided, and the annotators were informed that two out of three models had been post-edited, but none of them were informed which of the three systems had been post-edited. Evaluation measures are the classification accuracy of each annotator and the Fleiss’ K inter-annotator agreement: a low score on both can confirm that the models are indistinguishable.

4.2 Results and Discussion

LLMs leak Private Information Since LLMs tend to leak training data, we aim to quantify the amount of private information that can be retrieved from the pre-trained GPT-J. Unfortunately, GPT-J

makes no exception to the trend noticed by Huang et al. (2022) for the GPT-Neo models. In fact, also this model tends to generate PII.

In Table 2, it is possible to observe that Training Data Extraction Attacks that are based on Memorization are particularly effective against the GPT-J model: on average, the model tends to accurately predict the mail observed during training the 16% of the times.

It is worth noting the scale of the leakage: the model is originally prompted with 3238 examples. The column *Generated emails* reports the number of times during generation that the model answers with an email address, while *Leaked emails* reports the number of times the generation is correct, meaning that the generated email corresponds to the one observed in the training data. On average, 450.5 emails are correctly generated by those attacks: the privacy of a large number of people is threatened.

Moreover, as the attacker gets more information, the accuracy of the attacks gets higher: the accuracy of the attacks strongly depends on the length of the prompt. In fact, the lower accuracy – the number of correctly leaked email addresses over the total email addresses generated – that can be registered in Memorization Attacks is 12% : the model in that case is fed with a *context* prompt that is 50 tokens long. However, when the *context* prompt given to the model is composed of 200 tokens, the accuracy of the attack peaks at 18.2% with greedy decoding and 18% using beam search decoding.

The accuracy of the Association Attacks is much more modest. The results of those attacks against GPT-J model exhibit similar patterns to the one observed by Huang et al. (2022) against the GPT-Neo models. The larger number of email addresses leaked by those kind of attacks is 68, a modest number compared to the accuracy obtained in the

Memorization Attacks. However, in an adversarial scenario even low accuracy may cause harm to people. Hence, in the next Section, we will demonstrate the efficacy of PAE against both types of attacks.

PAE in batch editing Preserves Privacy In Table 2, it is possible to observe the reduced effectiveness of Memorization and Association attacks after the GPT-J model has undergone an editing process (Post-edit columns), with Post-edit results further divided into Implicit and Explicit categories. This Section investigates the impact of these edits, focusing on their efficacy against more or less informed attacks. We argue that PAE edits are effective if they can reduce the leakage of private information, regardless of the nature of the attack.

Post-edit results show, in fact, a significant reduction in the effectiveness of Association attacks. This reduction is particularly notable in scenarios where the number of leaked emails drops close to zero. For example, under the Implicit strategy, *zero-shot b* result in 0 leaked emails and 0 Attack accuracy, indicating a complete mitigation of the attack. However, some prompts still cause leakage; for instance, in prompt *zero-shot b* the edit reduces the number of email addresses leaked significantly but not completely (from 68 pre-edit to 48 post-edit in the Implicit). Crucially, while not perfect, the PAE edits – both Implicit and Explicit – always cause an increase in privacy protection, since reduce the number of email correctly leaked by Association Attacks.

However, it is crucial to consider the originally leaked emails when interpreting post-edit results. While a reduction to near-zero leakage is impressive, the impact is more pronounced when starting from a higher number of pre-edit leaks. For this reason, we focus on the discussion around the Memorization Attacks, that cause a larger number of private email addresses to be generated.

PAE is an effective solution against Memorization Attacks. In particular, the accuracy of the attacks steady decreases in each configuration. The average drop in accuracy after an Implicit edit is 5% and 4.5% after Explicit edit: this means that PAE is able to modify model parameter so that, on average, the 32% of the previously predicted email addresses are no more verbatim generated by the model using Implicit defense strategy, 29% with the explicit one. Against attacks with *context* prompt of 50 tokens, PAE effectiveness peaks, with

42.5% of the email addresses anonymized. As expected, more informed *context* prompts are more challenging: however, also with *context* prompts of 200 tokens, PAE make the accuracy attack drop from 0.18 to 0.12 in Greedy decoding and to 0.138 using Beam Search in the Implicit edit, and from 0.18 to 0.134 in greedy decoding and to 0.141 using Beam Search in the Explicit edit. In general, studying the effect of the decoding algorithm on the attacks accuracy we can state that this factor does not influence much the results: under Memorization Attacks only a slight difference in term of accuracy can be registered. From this analysis can conclude that PAE can help in protecting privacy.

Finally, it is possible to notice that there is a consistent difference between Implicit and Explicit post-edit results. Explicit edits generally result in a slightly higher number of leaked emails and attack accuracy compared to Implicit edits, especially under Memorization Attacks. For example, in the case of a *context* prompt of 200, the Explicit edit cause a larger number of emails to be correctly generated (396 in greedy decoding, 405 in beam search decoding) than the corresponding Implicit edit (368 email addresses leaked in greedy decoding, 394 in beam search decoding)

In summary, post-edit measures, particularly implicit edits, demonstrate a strong capability to safeguard email data from various attack strategies, significantly lowering both the number of leaked emails and the attack accuracy across different configurations. For our experiments, we adopted the "one model, n edits" philosophy. This approach is based on the practical scenario where an LLM producing private data needs to be edited for a large number of potentially threatened individuals: with PAE, the model owner can perform a single edit to the model parameters to reduce privacy risks. Our investigation on the emails that the model generates when subjected to Memorization and Association attacks confirm that the Memorization Attacks are able to recover a larger number of private information also after the editing. However, the evaluation confirms the effectiveness of our model editing techniques in preventing the disclosure of private data since also informed attacks –like the Memorization Attacks – are less effective on edited models. This setup is particularly challenging because it requires analyzing the impact of n modifications at the same time. In the next Section, we also demonstrate that the original language model is not negatively influenced by PAE.

Automatic Evaluation	BLEU	E_1-O	0.808(± 0.198)
		E_2-O	0.793(± 0.199)
		E_1-E_2	0.790(± 0.203)
	METEOR	E_1-O	0.841(± 0.173)
		E_2-O	0.826(± 0.172)
		E_1-E_2	0.824(± 0.183)
Manual Revision	Accuracy	0.35(± 0.07)	
	Fleiss' K	0.002	

Table 3: PAE preserves language model performances. In the Manual revision of the generation, annotators are not able to detect which model generated each paragraph. Moreover, on a larger scale of examples, the Automatic Evaluation reveals that the generation of *post-edit* LLMs E_1 and E_2 are both similar to one another (E_1-E_2) and with respect to the *pre-edit* LLM O (E_1-O and E_2-O)

PAE preserves the Language Modelling Capabilities PAE preserves privacy of people while not affecting the Language Modeling performances. Results of the evaluation can be found in the Table 3. The results of the *automatic assessment procedure* can quantitatively give us insight that models generations are, in fact, similar. Both according to BLEU metric and to METEOR, the systems generate (in greedy decoding) very similar paragraph when prompted with the same tokens. In particular, the post-edited models E_1 and E_2 – post-edited with implicit and explicit PAE – are similar to the original, pre-edited model O and are also similar with respect to each other. Finally, the *manual assessment procedure* suggest that the models are indistinguishable from one another. In fact, the annotators asked to detect which model is responsible for a generation among E_1 , E_2 , and O can only randomly guess, with an average accuracy on this classification task (0.35(± 0.07)) close to random choice. Also the very low agreement suggest that the three systems are indistinguishable.

This evaluation procedure can attest that the EAP is applicable because it not only helps to preserve user privacy, but also leaves the capabilities of the systems language model intact.

PAE is applicable with sequential editing Finally, in Table 3 it is possible to notice that the sequential update is definitely applicable with PAE. In this experiment, we perform sequential edit of the GPT-J model, varying the number of email addresses anonymized per edit, varying from 5 to 300. We indicate the number of anonymized emails per edit as batch size k : with $k < N$ we mimic the real world scenario of updating a model each time a privacy leak is detected. To understand whether

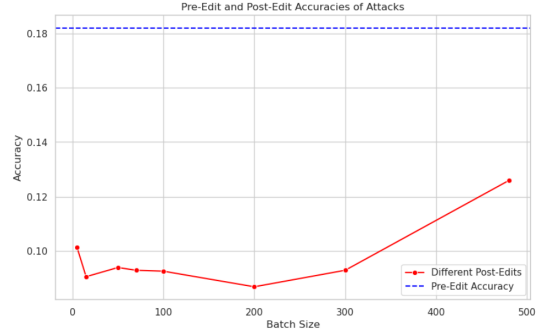


Figure 2: Memorization Attack against sequentially post-edited models. The smaller the batch size k , the larger the number of sequential updates are necessary to edit all the private emails addresses leaked by the original model. After the Sequential edits, the stronger Memorization Attack ($|p_M| = 200$) achieve similar performances at all the configurations.

sequential editing has a negative impact on effectiveness of the edit, we evaluate the effectiveness of PAE for each of the batch sizes in the Memorization Attack with the more effective of the prompts ($|p_M| = 200$). The results in Table 2 refers to a model post-edited with "implicit" PAE. As can be observed in Table 2, the accuracy of the edit is rather stable and similar to the results obtained in the batch editing scenario. Those results confirm the applicability of PAE also in sequential editing.

5 Conclusion

In this paper, we address the critical issue of private data leakage in Large Language Models (LLMs) due to their tendency to memorize training data. We propose Private Association Editing (PAE), a novel defense mechanism that effectively removes Personally Identifiable Information (PII) from LLMs without requiring retraining.

Our methodology involves a four-step procedure: detecting memorized PII, applying PAE cards, verifying resilience to targeted data extraction (TDE) attacks, and ensuring consistency in the post-edit LLMs. The PAE method stands out for its versatility and efficiency, allowing for small batch modifications and significantly enhancing the privacy of LLMs.

Our experiments demonstrate that the PAE approach is both effective and efficient in mitigating the risk of private data leakage. We believe PAE will be a valuable tool in the ongoing effort to protect data privacy in LLMs and encourage its adoption to prevent potential privacy violations as these models continue to be deployed in real-world applications.

705 Limitations

706 We outline some limitations and possible directions
707 for future research in enhancing data privacy in
708 Large Language Models (LLMs).

709 As the landscape of LLMs evolves, it may be
710 useful to extend the Private Association Editing
711 (PAE) mechanism to accommodate new types of
712 models and data. Currently, we apply our proposed
713 PAE method on a limited set of LLMs. A possible
714 extension could involve testing and refining
715 PAE across a broader spectrum of LLM architec-
716 tures and training datasets. Our approach focuses
717 on removing Personally Identifiable Information
718 (PII) from LLMs without retraining. However, this
719 method might not address all types of sensitive data.
720 Future research could explore additional techniques
721 to enhance the comprehensiveness of PII removal.
722 While PAE shows promise in its current form, its
723 real-world applicability and scalability need thor-
724 ough validation. By addressing these limitations,
725 future research can further solidify the role of PAE
726 in safeguarding data privacy in LLMs and ensure
727 its robustness and adaptability in various contexts.

728 References

729 Stella Biderman, Hailey Schoelkopf, Quentin Anthony,
730 Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mo-
731 hammad Aflah Khan, Shivanshu Purohit, USVSN Sai
732 Prashanth, Edward Raff, Aviya Skowron, Lintang
733 Sutawika, and Oskar van der Wal. 2023. *Pythia:
734 A suite for analyzing large language models across
735 training and scaling*. *Preprint*, arXiv:2304.01373.

736 Sid Black, Stella Biderman, Eric Hallahan, Quentin An-
737 thony, Leo Gao, Laurence Golding, Horace He, Con-
738 nor Leahy, Kyle McDonell, Jason Phang, Michael
739 Pieler, USVSN Sai Prashanth, Shivanshu Purohit,
740 Laria Reynolds, Jonathan Tow, Ben Wang, and
741 Samuel Weinbach. 2022. *Gpt-neox-20b: An open-
742 source autoregressive language model*. *Preprint*,
743 arXiv:2204.06745.

744 Nicola De Cao, Wilker Aziz, and Ivan Titov. 2021. *Edit-
745 ing factual knowledge in language models*. *Preprint*,
746 arXiv:2104.08164.

747 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,
748 Katherine Lee, Florian Tramèr, and Chiyuan Zhang.
749 2023. *Quantifying memorization across neural lan-
750 guage models*. *Preprint*, arXiv:2202.07646.

751 Nicholas Carlini, Florian Tramèr, Eric Wallace,
752 Matthew Jagielski, Ariel Herbert-Voss, Katherine
753 Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar
754 Erlingsson, et al. 2021. Extracting training data from
755 large language models. In *30th USENIX Security
756 Symposium (USENIX Security 21)*, pages 2633–2650.

Leo Gao, Stella Biderman, Sid Black, Laurence Gold- 757
ing, Travis Hoppe, Charles Foster, Jason Phang, 758
Horace He, Anish Thite, Noa Nabeshima, Shawn 759
Presser, and Connor Leahy. 2020. *The pile: An
800gb dataset of diverse text for language modeling*.
Preprint, arXiv:2101.00027. 760
761
762

Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, 763
Ruobing Xie, Julian McAuley, and Wayne Xin 764
Zhao. 2024. *Large language models are zero-
shot rankers for recommender systems*. *Preprint*,
arXiv:2305.08845. 765
766
767

Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. 768
2022. *Are large pre-trained language models leaking
your personal information?* In *Findings of the Asso-
ciation for Computational Linguistics: EMNLP 2022*,
pages 2038–2047, Abu Dhabi, United Arab Emirates.
Association for Computational Linguistics. 769
770
771
772
773

Bryan Klimt and Yiming Yang. 2004. The enron corpus:
A new dataset for email classification research. In
European conference on machine learning, pages
217–226. Springer. 774
775
776
777

Kevin Meng, David Bau, Alex Andonian, and Yonatan
Belinkov. 2023a. *Locating and editing factual associ-
ations in gpt*. *Preprint*, arXiv:2202.05262. 778
779
780

Kevin Meng, Arnab Sen Sharma, Alex Andonian,
Yonatan Belinkov, and David Bau. 2023b. *Mass-
editing memory in a transformer*. *Preprint*,
arXiv:2210.07229. 781
782
783
784

Eric Mitchell, Charles Lin, Antoine Bosselut, Christo- 785
pher D. Manning, and Chelsea Finn. 2022. *Memory-based model editing at scale*. *Preprint*,
arXiv:2206.06520. 786
787
788

Yongyu Mu, Abudurexiti Reheman, Zhiqian Cao, 789
Yuchun Fan, Bei Li, Yinqiao Li, Tong Xiao, Chun-
liang Zhang, and Jingbo Zhu. 2023. *Augmenting
large language model translators via translation mem-
ories*. In *Findings of the Association for Computa-
tional Linguistics: ACL 2023*, pages 10287–10299,
Toronto, Canada. Association for Computational Lin-
guistics. 790
791
792
793
794
795
796

Milad Nasr, Nicholas Carlini, Jonathan Hayase, 797
Matthew Jagielski, A Feder Cooper, Daphne Ippolito,
Christopher A Choquette-Choo, Eric Wallace, Flo-
rian Tramèr, and Katherine Lee. 2023. Scalable ex-
traction of training data from (production) language
models. *arXiv preprint arXiv:2311.17035*. 798
799
800
801
802

Vaidehi Patil, Peter Hase, and Mohit Bansal. 2023. 803
*Can sensitive information be deleted from llms?
objectives for defending against extraction attacks*.
Preprint, arXiv:2309.17410. 804
805
806

Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim 807
Rocktäschel, Yuxiang Wu, Alexander H. Miller,
and Sebastian Riedel. 2020. *How context af-
fects language models’ factual predictions*. *ArXiv*,
abs/2005.04611. 808
809
810
811

812	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	Problems, methods, and opportunities. <i>Preprint</i> ,	871
813	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	arXiv:2305.13172.	872
814	Alexander Miller. 2019. Language models as knowl-		
815	edge bases? In <i>Proceedings of the 2019 Confer-</i>	6 Appendix	873
816	<i>ence on Empirical Methods in Natural Language Pro-</i>		
817	<i>cessing and the 9th International Joint Conference</i>	6.1 All the prompts of the manual assessment	874
818	<i>on Natural Language Processing (EMNLP-IJCNLP)</i> ,	procedure	875
819	pages 2463–2473, Hong Kong, China. Association	• "My name is";	876
820	for Computational Linguistics.		
821	Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie	• "The weather today is";	877
822	Millican, Jordan Hoffmann, Francis Song, John		
823	Aslanides, Sarah Henderson, Roman Ring, Susan-	• "Ever since I was a child, I've always liked";	878
824	nah Young, Eliza Rutherford, Tom Hennigan, Ja-		
825	cob Menick, Albin Cassirer, Richard Powell, George	• "My dear friend Mary";	879
826	van den Driessche, Lisa Anne Hendricks, Mari-		
827	beth Rauh, Po-Sen Huang, Amelia Glaese, Jo-	• "Swimmers are usually";	880
828	hannes Welbl, Sumanth Dathathri, Saffron Huang,		
829	Jonathan Uesato, John Mellor, Irina Higgins, Anto-	• "Modern art is";	881
830	nia Creswell, Nat McAleese, Amy Wu, Erich Elsen,		
831	Siddhant Jayakumar, Elena Buchatskaya, David Bud-	• "The Industrial Revolution";	882
832	den, Esme Sutherland, Karen Simonyan, Michela Pa-		
833	ganini, Laurent Sifre, Lena Martens, Xiang Lorraine	• "Follow those steps to cook";	883
834	Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena		
835	Gribovskaya, Domenic Donato, Angeliki Lazaridou,	• "It is forbidden to";	884
836	Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsim-		
837	poukelli, Nikolai Grigorev, Doug Fritz, Thibault Sot-	• "It is very likely".	885
838	tiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong,		
839	Daniel Toyama, Cyprien de Masson d'Autume, Yujia		
840	Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin,		
841	Aidan Clark, Diego de Las Casas, Aurelia Guy,		
842	Chris Jones, James Bradbury, Matthew Johnson,		
843	Blake Hechtman, Laura Weidinger, Iason Gabriel,		
844	William Isaac, Ed Lockhart, Simon Osindero, Laura		
845	Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub,		
846	Jeff Stanway, Lorraine Bennett, Demis Hassabis, Ko-		
847	ray Kavukcuoglu, and Geoffrey Irving. 2022. Scaling		
848	language models: Methods, analysis insights from		
849	training gopher. <i>Preprint</i> , arXiv:2112.11446.		
850	Federico Ranaldi, Elena Sofia Ruzzetti, Dario Ono-		
851	rati, Leonardo Ranaldi, Cristina Giannone, Andrea		
852	Favalli, Raniero Romagnoli, and Fabio Massimo Zan-		
853	zotto. 2024. Investigating the impact of data con-		
854	tamination of large language models in text-to-sql		
855	translation. <i>Preprint</i> , arXiv:2402.08100.		
856	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-		
857	6B: A 6 Billion Parameter Autoregressive Lan-		
858	guage Model. https://github.com/kingoflolz/		
859	mesh-transformer-jax .		
860	Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie,		
861	Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi		
862	Xiong, Ritik Dutta, Rylan Schaeffer, Sang T. Truong,		
863	Simran Arora, Mantas Mazeika, Dan Hendrycks, Zi-		
864	nan Lin, Yu Cheng, Sanmi Koyejo, Dawn Song, and		
865	Bo Li. 2024. Decodingtrust: A comprehensive as-		
866	sessment of trustworthiness in gpt models. <i>Preprint</i> ,		
867	arXiv:2306.11698.		
868	Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng,		
869	Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu		
870	Zhang. 2023. Editing large language models:		