## Soft Token Attacks Cannot Reliably Audit Unlearning in Large Language Models

Anonymous ACL submission

#### Abstract

Large language models (LLMs) have become increasingly popular. Their emergent capabilities can be attributed to their massive training datasets. However, these datasets often contain undesirable or inappropriate content, e.g., harmful texts, personal information, and copyrighted material. This has promoted research into *machine unlearning* that aims to remove information from trained models. In particular, *approximate unlearning* seeks to achieve information removal by strategically editing the model rather than complete model retraining.

002

012

017

023

024

027

033

034

Recent work has shown that soft token attacks (STA) can successfully extract purportedly unlearned information from LLMs, thereby exposing limitations in current unlearning methodologies. In this work, we reveal that STAs are an inadequate tool for auditing unlearning. Through systematic evaluation on common unlearning benchmarks (Who Is Harry Potter? and TOFU), we demonstrate that such attacks can elicit any information from the LLM, regardless of (1) the deployed unlearning algorithm, and (2) whether the queried content was originally present in the training corpus. Furthermore, we show that STA with just a few soft tokens (1 - 10) can elicit random strings over 400-characters long. Thus showing that STAs are too powerful, and misrepresent the effectiveness of the unlearning methods.

> Our work highlights the need for better evaluation baselines, and more appropriate auditing tools for assessing the effectiveness of unlearning in LLMs.

#### 1 Introduction

In recent years, large language models (LLMs)
have undergone substantial advancements, leading
to enhanced performance and widespread adoption. LLMs have demonstrated exceptional performance in various downstream tasks, such as machine translation (Zhu et al., 2023), content genera-

tion (Acharya et al., 2023), and complex problemsolving (Xi et al., 2025). Their performance is attributed to their large-scale architectures that require datasets consisting of up to billions of tokens to train effectively (Kaplan et al., 2020). These datasets are typically derived from large-scale corpora sourced from public internet text. However, such datasets inadvertently contain harmful or inappropriate content, including instructions for hazardous activities (e.g., bomb-making), violent or explicit material, private information, or copyrighted content unsuitable for applications. Given the sensitive nature of such data, it may be necessary to remove it from the LLM to comply with the local regulations, or internal company policies. 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

Machine unlearning is a tool for removing information from models (Cao and Yang, 2015; Bourtoule et al., 2021a). Approximate unlearning usually refers to removing information from models without resorting to retraining them from scratch (Zhang et al., 2024a; Eldan and Russinovich, 2023a; Izzo et al., 2021), ensuring that the resulting model deviates from a fully retrained version within a bounded error. While numerous studies have proposed various unlearning algorithms, most lack formal guarantees of effectiveness. In fact, prior research has demonstrated that many unlearning techniques can be circumvented through simple rephrasings of the original data (Shi et al., 2024). Recent work has shown that a *soft token* attack (STA) can be used to elicit harmful completions and extract supposedly unlearned information from models (Schwinn et al., 2024; Zou et al., 2024).

In this work, we introduce a simple framework for *auditing unlearning* and demonstrate that *STA*s are overly powerful, and hence, inappropriate for verifying the effectiveness of approximate unlearning. We show that the auditor can elicit any information from the model, regardless of its training data. Our work highlights the need for better un-

- 084 085 086
- 08
- 80
- 09
- 0
- 0
- 090

0

099

1

101 102

103 104

105

106

107

108 109

# 110

## 111

112

113

114

115

116

117

118

119

## 2.1 Background

2

in Section 7.

**Large language models** (LLMs) process input text through an auto-regressive framework. Given an input sequence of tokens  $x_{1:t}$ , the model computes the conditional probability distribution  $p(x_{t+1}|x_{1:t})$ over the vocabulary at each time-step. The likelihood of the sequence is given by:

**Background & Related work** 

learning auditing baselines and methodologies.

1. We introduce a simple auditing framework for

2. We show that STAs effectively elicit un-

learned information in all tested unlearning

methods and benchmark datasets (Who Is

Harry Potter?, and TOFU). Additionally, we

show that STAs also elicit information in the

base models that were not fine-tuned on the

3. We further demonstrate that the *STA*'s are inappropriate for evaluating unlearning – we

show that a single soft token can elicit 150

random tokens, and ten soft tokens can elicit

The remainder of this paper is organized as fol-

lows: In Section 2 we provide an overview of the

necessary background, and the related work. Sec-

tion 3 introduces a general auditing framework, and

instantiates it using STA. In Section 4, we demon-

strate the efficacy of STA, and subsequently its

failure, as a tool for auditing unlearning in LLMs.

In Section 5 we discuss additional considerations

for auditing unlearning in LLMs. We conclude the

paper in Section 6, and highlight some limitations

over 400 random tokens (Section 4.3).

benchmark datasets (Section 4.2).

We claim the following contributions:

unlearning in LLMs (Section 3.2).

$$\log p(x_{t+1}|x_{1:t}) = \sum_{t=1}^{T} \log p(x_t|x_{1:t-1})$$
 (1)

120At inference time, the tokens is generated itera-121tively by sampling the next token  $x_{t+1}$  from this122distribution (e.g., via greedy decoding or nucleus123sampling (Holtzman et al., 2019)), then appending124it to the context  $x_{1:t}$  for the subsequent step.

125Machine unlearning is a tool for removing infor-<br/>mation from models. Consider a machine learning<br/>model f optimized over a training dataset  $D_{train}$ .128When a data owner submits an unlearning request

to remove a specified subset  $D_{forget} \in D_{train}$ , the objective of machine unlearning is to produce an unlearned model  $f_u$  that eliminates the influence of  $D_{forget}$ . Machine unlearning methodologies are categorized into two paradigms – exact, and approximate unlearning.

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

169

170

171

172

173

174

175

176

**Exact unlearning** ensures the output distribution of  $f_u$  is statistically indistinguishable from that of a model retrained exclusively on the retained dataset  $D_{retain} = D_{train}/D_{forget}$ . This guarantees provable data removal, satisfying:

$$p(f_u(x) = y) = p(f_{ret}(x) = y)$$
  
s.t.  $\forall (x, y) \in D_{train},$  (2)

where  $f_{ret}$  denotes a model retrained from scratch on  $D_{retain}$  – which is the most straightforward way of achieving exact unlearning.

The process can be made more efficient by splitting the  $D_{train}$  into overlapping chunks, and training an ensemble of models (Bourtoule et al., 2021b). During an unlearning request, only the models containing the requested records need to be retrained. For certain classes of models, it is possible to achieve exact unlearning without retraining, e.g. ECO adapts the Cauwenberghs and Poggio (CP) algorithm for exact unlearning within LeNet (Huang et al.), and MUSE relabels the target data to achieve unlearning for over-parameterized linear models (Yang et al., 2024).

**Approximate unlearning**, sometimes called *in*exact unlearning, relaxes the strict equivalence requirement, instead only requiring that  $f_u$  approximates  $f_{ret}$  within some bounded error. This paradigm relies on empirical metrics or probabilistic frameworks. In LLMs, approximate unlearning is typically accomplished by overwriting the information in the model (Eldan and Russinovich, 2023a; Wang et al., 2024), guiding the model away from it (Feng et al., 2024), or editing the weights and/or activations (Liu et al., 2024; Bhaila et al., 2024; Li et al., 2024; Tamirisa et al., 2024; Huu-Tien et al., 2024; Ashuach et al., 2024; Meng et al., 2022a,b).

## 2.2 Related work

While advances have been made in developing machine unlearning algorithms for LLMs, rigorous methodologies for auditing the efficacy of unlearning remain understudied. Adversarial soft token attacks (*STA*s) (Schwinn et al., 2024) and 5-shot in-context prompting (Doshi and Stickland, 2024)



Figure 1: Overview of the auditing process using  $A_{STA}$ . For a perfect unlearning method,  $A_o$  always correctly audits the model. On the other hand,  $A_{STA}$  can elicit the completion regardless of the information in the model – the audit is ineffective.

have been shown to recover unlearned knowledge in models. When model weights can be modified, techniques such as model quantization (Zhang et al., 2024e) and retraining on a partially unlearned dataset (Łucki et al., 2024; Hu et al., 2024) have also proven effective in recalling forgotten information. (Lynch et al., 2024) examined eight methods for evaluating LLM unlearning techniques and found that their latent representations remained similar. News and book datasets are used to analyze unlearning algorithms from six different perspectives (Shi et al., 2024). It was shown that finetuning on unrelated data can restore information unlearned from the LLM (Qi et al., 2024), indicating the existing unlearning methods do not actually remove the information but learn a refusal filter instead. Several benchmarks have been developed to evaluate the existing unlearning algorithms. Besides, an unlearning benchmark was introduced based on fictitious author information (Maini et al., 2024a). For real-world knowledge unlearning, Real-World Knowledge Unlearning (RWKU) used 200 famous people as unlearning targets (Jin et al., 2024), while WDMP focused on unlearning hazardous knowledge in biosecurity, cybersecurity (Li et al., 2024).

177

178

182

183

190

191

194

195

198

199

200

203

207

## **3** Auditing with Soft Token Attacks

## 3.1 Adversarial prompts

An adversarial prompt  $x_a$ , is an input prompt to the LLM, obtained by applying the transform  $T(\cdot)$ to the base prompt  $x_p$ ,  $x_a = T(x_p, aux)$  in order to elicit a desired completion c. T can be any function that swaps, removes or adds tokens; *aux* denotes any additional needed information. However, such arbitrary attacks are expensive to optimize<sup>1</sup>, and difficult to reason about. In practice, T optimizes an *adversarial suffix*  $x_s$  that is appended to  $x_p$  to elicit c (Zou et al., 2023). Specifically, we optimize the probability:

$$Prob = P(c|x_p \oplus x_s). \tag{3}$$

209

210

211

213

214

215

216

217

218

220

221

223

224

225

226

227

228

229

230

231

232

An adversary with white-box access to the LLM, can instead mount the attack in the *embedding space* i.e. modify the *soft tokens*:

$$Prob = P(c|embed(x_p) \oplus embed(x_s)). \quad (4)$$

In this case, T uses the gradient from the LLM to update  $x_s$ .

#### 3.2 Unlearning auditor

An *oracle* auditor  $A_o$  takes an unlearned model  $f_u$ and the candidate sentences  $x_c \in X_c$  and outputs a ground truth, binary decision  $a = \{0, 1\}$  indicating whether the given records was part of  $D_{train}$  of:

$$a = A_o(f_u, X_c = D_{forget}, aux)$$
(5)

 $A_o$  is unrealistic in many scenarios; however, it can be easily instantiated for exact unlearning where  $A_o$  knows the training data associated with  $f: aux = \{D_{retain}\}.$ 

<sup>&</sup>lt;sup>1</sup>Suffix-only attacks allow efficient use of the KV-cache.

On the other hand, a realistic unlearning  $A_u$  takes an  $f_u$ , and  $D_{forget}$  and outputs a score s = (0, 1) indicating whether the records were in  $D_{train}$ :

$$s = A_u(f_u, \{x_c\} = D_{forget}, aux = \emptyset) \quad (6)$$

 $A_u$  represents cases where users remove information from models that they did not create, e.g. to prevent harmful outputs.

In this work, we instantiate the soft token attack auditor  $A_{STA}$  based on the soft token attacks (STAs) against unlearning (Schwinn et al., 2024; Zou et al., 2024). Our auditor compares the relative difficulty of eliciting c for  $f_{ft}$  and  $f_u$ . The unlearning procedure is effective if eliciting completions using  $f_u$  is more difficult than  $f_{ft}$ .

$$s = A_{STA}(f_u, \{x_c\} = D_{forget}, aux = \{f_{ft}\}).$$
(7)

Figure 1 gives a complete overview of the auditing procedure, and the difference between  $A_0$  and  $A_{STA}$ . In Table 1, we summarize the notation.

In the next Section, we show that  $A_{STA}$  cannot reliably audit LLMs.

#### 4 Evaluation

234

237

238

240

241

243

245

247

248

249

251

259

260

261

263

271

272

273

276

277

#### 4.1 Experiment setup

**Datasets.** To attack unlearning and evaluate its effectiveness, we use two popular benchmark datasets: 1) *Who Is Harry Potter?* (Eldan and Russinovich, 2023a), a benchmark that intends to remove information about the world of Harry Potter and the associated characters; *WHP* hereafter. 2) *TOFU* (Maini et al., 2024a), a dataset of fictional writers that are guaranteed to be absent in the LLM's training data<sup>2</sup>.

WHP does not publish a complete dataset. For that reason we use the passages included in the associated Hugging Face page (Eldan and Russinovich, 2023b). Additionally, we augment it with  $20 (x_p \rightarrow c)$  pairs generated with Llama2-7b-chathf. These contain general trivia about the Harry Potter universe.

For *TOFU*, we use the 10% forget to 90% retain split provided by the authors (Maini et al., 2024b). **Models & environment.** For all experiments, we use Llama-2-7b-chat-hf (Touvron et al., 2023) (Llama2), and Llama-3-8b-instruct (Meta, 2024) (Llama3) downloaded from Hugging Face. We

STA	soft token attack
$A_o$	oracle auditor
$A_{STA}$	STA auditor
$x_p$	base prompt (benign)
$x_s$	adversarial suffix
$x_a$	adversarial prompt ( $x_p \oplus x_s$ )
c	target completion
$f_{\emptyset}$	base model
$f_{ft}$	fine-tuned model
$f_u$	unlearned model
$f_{u-*}$	model unlearned using *
$D_{train}$	training data
$D_{forget}$	forget data
$D_{retain}$	retain data

Table 1: Summary of the notation. '\*' is replaced with the specific unlearning method.

get the unlearned *WHP* model from it's Hugging Face repository (Eldan and Russinovich, 2023b) (Llama2-*WHP*).

We implement *STA* using LLMart (Cornelius et al., 2025) – a PyTorch and Hugging Face-based library for crafting adversarial prompts. We use implementations of the unlearning methods from the *TOFU* (Maini et al., 2024c), and NPO (Zhang et al., 2024b) repositories. We benchmark the attack against seven different unlearning algorithms: gradient ascent (GA), gradient difference (GDF) (Liu et al., 2022), refusal (IDK) (Rafailov et al., 2024), knowledge distillation (KL) (Hinton, 2015), negative preference optimization (NPO) (Zhang et al., 2024c), NPO-GDF, NPO-KL.

We run our experiments on a machine equipped with Intel Xeon Gold 5218 CPU, eight NVIDIA A6000, and 256 GB of RAM.

#### 4.2 Auditing with attacks

Who Is Harry Potter?. To elicit information about the Harry Potter universe, we initialize the soft tokens using randomly selected hard tokens, and append them to the prompt  $(embed(x_a = x_p \oplus x_s))$ . We then train the soft prompt using AdamW (Loshchilov and Hutter, 2019) for up to 3000 iterations; using learning 0.005, and default  $\beta s$  – we reiterate that the  $x_p$  does not change, only the embedded suffix does. If the optimization fails, we double the the number of soft tokens up to the maximum of 16. We rerun the experiment five times, and report the mean and standard deviation across all prompt and reruns. In Table 2 we report the average number of soft tokens needed to elicit a

<sup>&</sup>lt;sup>2</sup>This cannot be guaranteed for models published after *TOFU*.

completion. *WHP* \* denotes the unlearned model with different prompt templates.

311

312

313

314

316

317

319

321

322

323

324

325

327

329

333

334

335

337

339

341

347

348

351

354

362

Our results show that all information can be generated with  $\approx 4-6$  added soft tokens. For all pairs of models, we conduct a *t-test* under the null hypothesis  $\mathcal{H}_0$  of *equivalent population distributions* with  $\alpha = 0.05$ . We use an unpaired Welch's ttest since sample variances are not equal (WELCH, 1947). We cannot reject the hypothesis for any of the pairs i.e. p > 0.05. In other words, for all models, there is not enough evidence to say that eliciting completions is more difficult.

Additionally, we observe that the ease of eliciting the completions changes depending on the prompt template. We conducted our initial exploratory experiments in the chat setting with a chat prompt template. We notice that the model (WHP +chat) reveals all unlearned information with manually paraphrased prompts. Furthermore, when using the example prompts in the corresponding Hugging Face repository (Eldan and Russinovich, 2023b), the model (WHP) would often begin the response with a double new line (\n\n). We suspect that the provided unlearned model is overfit to "prompt\n\n completion". To run our evaluation in the most favorable setting, we report all three. Our results show that attacking is the easiest for WHP +chat, and the most difficult for WHP +\n\n. Given these discrepancies, and the lack of a standard WHP dataset, we believe it is not a good unlearning benchmark, despite its popularity.

Also, in our dataset there are three challenging outlier prompts that require 16 soft tokens, unlike other prompts. Filtering these out results in 4.05, 4.60, 1.61 average required tokens for WHP, WHP +\n\n, and WHP +chat respectively.

**TOFU.** For *TOFU*, we follow the same setup as for *WHP*– we initialize soft tokens using random hard tokens, and append them; we then train the soft prompt using AdamW for up to 3000 iterations; using learning 0.005, and  $\beta s = (0.9, 0.999)$ ; we double the soft tokens if the optimization fails; we rerun the experiments five times and report the averages across all prompts. In Table 3, we report the number of soft tokens required to elicit the completions.  $f_{\emptyset}$  refers to the unmodified baseline model,  $f_{ft}$  corresponds to the models fine-tuned on *TOFU*, followed by the unlearned models.

For all methods, we can elicit the completions with  $\approx 3$  appended soft tokens. Similarly to *WHP*, for all possible pairs of models (within the same architecture), we conduct a *t-test* under the null hy-

Prompt	Model		
template	Llama2-WHP	Llama3	
N/A	N/A	$5.61 \pm 6.32$	
WHP	$4.63 \pm 3.69$	N/A	
WHP + n n	$6.50\pm5.13$	N/A	
WHP +chat	$4.12\pm5.53$	N/A	

Table 2: Number of soft tokens required to elicit a completion for a fixed number of iterations. Soft tokens are appended to the prompt. Results are averaged over all prompts in the *WHP* set and over five runs for each prompt. When we increase the maximum iterations to 10,000 we can elicit **all** completions with 1 - 2 soft tokens. We do not report the results for Llama2 hf because it was used to generated the data. For comparison we also report the results for Llama3 without any prompt template (N/A).

Unloaming mathed	Model		
Omeaning method	Llama2	Llama3	
$f_{\emptyset}$ (none)	$3.07 \pm 3.25$	$3.11 \pm 3.15$	
$f_{ft}$ (none)	$2.95 \pm 3.35$	$3.21 {\pm} 3.19$	
$f_{u-IDK}$	$3.40 \pm 3.20$	$3.33 {\pm} 3.09$	
$f_{u-GA}$	$3.34{\pm}3.97$	$3.21 {\pm} 3.87$	
$f_{u-GDF}$	$3.06 \pm 3.34$	$3.11 \pm 3.40$	
$f_{u-KL}$	$3.08 \pm 3.31$	$3.12 {\pm} 3.17$	
$f_{u-NPO}$	$3.11 \pm 3.27$	$3.12 {\pm} 3.27$	
$f_{u-NPO-GDF}$	$3.15 \pm 3.24$	$3.16 {\pm} 3.16$	
$f_{u-NPO-KL}$	$3.23 {\pm} 3.62$	$3.24 {\pm} 3.57$	

Table 3: Number of soft tokens required to elicit a completion for a fixed number of iterations. Soft tokens are appended to the prompt. Results are averaged over all prompts in the *TOFU* set and over five runs for each prompt. When we increase the maximum iterations to 10,000 we can elicit **all** completions with 1 - 2 soft tokens.

pothesis  $\mathcal{H}_0$  of *equivalent population distributions* with  $\alpha = 0.05$ . We use an unpaired Welch's t-test since sample variances are not equal.

We cannot reject the hypothesis for any of the pairs i.e. p > 0.05;  $f_{u-IDK}$  vs  $f_{ft}$  (for Llama2) gives the lowest p-value of 0.509. In other words, for all models (regardless if trained on *TOFU*, or used unlearning method), there is not enough evidence to say that eliciting completions is more difficult.

One could argue that used unlearning methods are not effective (when comparing  $f_{ft}$  vs  $f_{u-*}$ ), hence they require similar numbers of soft tokens.<sup>3</sup>

363

<sup>&</sup>lt;sup>3</sup>Most of these approaches have been in fact been shown ineffective, and susceptible to simple paraphrasing.

However, the same holds when compared to  $f_{\emptyset}$ . In the next section, we demonstrate that the result cannot be attributed to the (in-)effectiveness of the unlearning methods but rather the power of *STA*.

## 4.3 Eliciting random strings

381

383

387

394

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

494

425

To further show the power of *STA*s, we use them to elicit *random strings*. Unlike natural text, the chance of a random string appearing in the training set is negligible. Also, preceding tokens do not inform the selection of the next token. In the following experiments, we pick characters uniformly at random in the range 33-126 of the ASCII table (asciitable.com).

To elicit a random string, we initialize the soft prompt using randomly selected tokens. Unlike in the WHP and TOFU experiments, there is only the soft prompt. We then train the soft prompt using AdamW for up to 3000 iterations per soft token; using learning rate 0.005, and  $\beta s = (0.9, 0.999)$ .

In Figure 2, we report the longest elicited string for a given number of soft tokens. We repeat the experiment five times – e.g., the first marker implies that for each of the five tested random strings of length 150, we found an effective soft prompt. We observe that not all initializations and seed configurations succeed, in which case a run needs to be restarted with a different seed. If the loss plateaus around 25% of the iterations, we restart the run. However, no single string was restarted more than ten times. We experimented with learning rate schedulers but they did not improve the search.

Our results show that *STA*s can be used to elicit completely random strings, thus undermining their application for auditing unlearning. Due to limited computational resources and long run-times, our evaluation is limited to 10 soft tokens, and 400-character long random strings. This does not provide a bound on the longest string that can be elicited.

Next, we aim to answer why eliciting strings is possible. Prompt-tuning (Lester et al., 2021) is a *performance efficient fine-tuning technique* in which instead of training all weights, one trains only a soft prompt added to the input. *STA* s can be viewed as an extreme case of prompt-tuning, where instead of training over many prompts, one trains an attack per each prompt. Hence, an LLM that outputs a completion that it was trained on is an expected behavior. We urge against misinterpreting the results and declaring techniques ineffective.



Figure 2: Numbers of random characters generated for the given soft prompt length. A single soft token can force over 150 random characters – more than any text in the *TOFU* benchmark. With 10 soft tokens, it is possible to generate over 400 random characters.

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

### 5 Discussion & Conclusion

Auditing with hard prompts. Attacks such as greedy coordinate gradient (Zou et al., 2023) optimize the attack prompt in the *hard* token space instead of the soft token space. Hence, they are weaker at eliciting completions. On one hand, this might make them more suitable for auditing unlearning. On the other hand, due to their computational requirements, they are often used to force only the beginning of a harmful completion (e.g. *Sure, here's how to build a bomb...*) with the hope that the LLM follows. It is unclear whether this would be sufficient to produce specific unlearned passages. We see it as an interesting direction for future work.

**Unlearning vs jail-breaking.** Our findings are applicable to the jail-breaking community as well. Prior work (Zhang et al., 2024d) hinted that unlearning and preventing harmful outputs can be viewed as the same task – removing or suppressing particular information. *STA*'s and fine-tuning attacks (Hu et al., 2024) are useful tools for evaluating LLMs in powerful threat models. It was shown that fine-tuning on benign data, or data unrelated to the unlearning records (for jail-breaking and unlearning respectively) can restore undesirable behavior (Łucki et al., 2024).

**Variation in gradient-based learning.** Prior work showed that removing training records from the training set, and repeating the training can result in the same final model (Thudi et al., 2022) depending on the random seed. Even though a record was

6

507

part of the training run, its influence might be minimal, making unlearning unnecessary. Similarly, it was shown that SGD has intrinsic privacy guarantees, assuming there exists a group of similar records (Hyland and Tople, 2022). Thus, algorithmic auditing of unlearning might not be possible, and one would have to rely on verified or attested procedures instead (Eisenhofer et al., 2023), regardless of their impact on the model.

458

459

460

461

462

463

464

465

466

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

Distinguishing learned soft tokens. Even though, 467 in most our results, the number of soft tokens re-468 quired to elicit a completion is the same, we at-469 tempted to distinguish between them. To this end, 470 471 we take all single-token STAs optimized for TOFU (Table 3) and assign a label  $y = \{0, 1\}$ : y = 0472 for  $f_{\emptyset}$ , and y = 1 for  $f_{ft}$  and the unlearned mod-473 els. We then train a binary classifier using  $f_{\emptyset}$  and 474  $f_{ft}$ . While we are able to overfit it and distinguish 475 between  $f_{\emptyset}$  and  $f_{ft}$ , we were not able to train a 476 model that would generalize to the unlearned mod-477 els, and decisively assign a class. Our approach 478 is similar to Dataset Inference (Maini et al., 2021, 479 2024d) which showed there can be distributional 480 differences between the models, depending on the 481 data they were trained on. Further investigation 482 into *what* soft tokens are learned during the audit 483 is an interesting direction for future work. 484

## 6 Conclusion

In this work, we show that soft token attacks (*STA*s) cannot reliably distinguish between base, fine-tuned, and unlearned models. In all cases, the auditor can elicit all unlearned information by appending optimized soft prompts to the base prompt. Additionally, we show that *STA* with a single soft token can elicit 150 random characters, and over 400 with soft tokens.

Our work demonstrates that machine unlearning in LLMs needs better evaluation frameworks. While many unlearning methods can be broken by simple paraphrasing of original prompts, or by finetuning on partial unlearned data or even *unrelated data*, *STA* misrepresents their efficacy.

## 7 Limitations & ethical considerations

Limitations. Due to computational constraints our
work is limited to 7-8 billion parameter models.
Nevertheless, given that LLMs' expressive power
increases with size (Kaplan et al., 2020), our results should hold for larger LLMs. Our evaluation
with random strings could be extended to verify if

there is a clear and generalizable dependency between the number of soft tokens and the maximum number of generated characters.

**Ethical considerations.** In this work, we show that an auditor (a user) with white-box access to the model, and sufficient compute can elicit any text from the LLM. While it does require knowing the target completion for a given prompt, it is likely that partial completions might be enough, thus allowing the user to elicit harmful information. This may be particularly dangerous in settings where the user has approximate knowledge of the information that had been scrubbed off the LLM.

#### References

Arkadeep Acharya, Brijraj Singh, and Naoyuki Onoe. 2023. Llm based generation of item-description for recommendation system. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 1204–1207.

asciitable.com. Ascii table.

- Tomer Ashuach, Martin Tutek, and Yonatan Belinkov. 2024. Revs: Unlearning sensitive information in language models via rank editing in the vocabulary space. *arXiv preprint arXiv:2406.09325*.
- Karuna Bhaila, Minh-Hao Van, and Xintao Wu. 2024. Soft prompting for unlearning in large language models. *arXiv preprint arXiv:2406.12038*.
- Lucas Bourtoule, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. 2021a. Machine unlearning. In 2021 IEEE Symposium on Security and Privacy (SP), pages 141–159. IEEE.
- Lucas Bourtoule, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. 2021b. Machine unlearning. In 2021 IEEE Symposium on Security and Privacy (SP), pages 141–159. IEEE.
- Yinzhi Cao and Junfeng Yang. 2015. Towards making systems forget with machine unlearning. In 2015 *IEEE symposium on security and privacy*, pages 463– 480. IEEE.
- Cory Cornelius, Marius Arvinte, Sebastian Szyller, Weilin Xu, and Nageen Himayat. 2025. LLMart: Large Language Model adversarial robutness toolbox.
- Jai Doshi and Asa Cooper Stickland. 2024. Does unlearning truly unlearn? a black box evaluation of llm unlearning methods. *Preprint*, arXiv:2411.12103.
- Thorsten Eisenhofer, Doreen Riepel, Varun Chandrasekaran, Esha Ghosh, Olga Ohrimenko, and Nicolas Papernot. 2023. Verifiable and provably secure machine unlearning. *Preprint*, arXiv:2210.09126.

- 559 560 562 565 567 568 569 570 571 574 575 577 578 581 585 586 587 588 590 594 597 599
- 605 606 607
- 610
- 611 612
- 613

- Ronen Eldan and Mark Russinovich. 2023a. Who's harry potter? approximate unlearning in llms. arXiv preprint arXiv:2310.02238.
- Ronen Eldan and Mark Russinovich. 2023b. Who's harry potter? approximate unlearning in llms.
- XiaoHua Feng, Chaochao Chen, Yuyuan Li, and Zibin Lin. 2024. Fine-grained pluggable gradient ascent for knowledge unlearning in language models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 10141-10155.
- Geoffrey Hinton. 2015. Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. arXiv preprint arXiv:1904.09751.
- Shengyuan Hu, Yiwei Fu, Steven Wu, and Virginia Smith. 2024. Jogging the memory of unlearned models through targeted relearning attacks. In ICML 2024 Workshop on Foundation Models in the Wild.
- Yu-Ting Huang, Pei-Yuan Wu, and Chuan-Ju Wang. Eco: Efficient computational optimization for exact machine unlearning in deep neural networks. In 2nd Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ ICML 2024).
- Dang Huu-Tien, Trung-Tin Pham, Hoang Thanh-Tung, and Naoya Inoue. 2024. On effects of steering latent representation for large language model unlearning. arXiv preprint arXiv:2408.06223.
- Stephanie L. Hyland and Shruti Tople. 2022. An empirical study on the intrinsic privacy of sgd. Preprint, arXiv:1912.02919.
- Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. 2021. Approximate data deletion from machine learning models. In International Conference on Artificial Intelligence and Statistics, pages 2008-2016. PMLR.
- Zhuoran Jin, Pengfei Cao, Chenhao Wang, Zhitao He, Hongbang Yuan, Jiachun Li, Yubo Chen, Kang Liu, and Jun Zhao. 2024. Rwku: Benchmarking realworld knowledge unlearning for large language models. Preprint, arXiv:2406.10890.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. Preprint, arXiv:2001.08361.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 3045-3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Nathaniel Li, Alexander Pan, Anjali Gopal, Summer 614 Yue, Daniel Berrios, Alice Gatti, Justin D Li, Ann-615 Kathrin Dombrowski, Shashwat Goel, Long Phan, 616 et al. 2024. The wmdp benchmark: Measuring and re-617 ducing malicious use with unlearning. arXiv preprint 618 arXiv:2403.03218. 619 Bo Liu, Qiang Liu, and Peter Stone. 2022. Continual 620 learning and private unlearning. In Conference on 621 Lifelong Learning Agents, pages 243–254. PMLR. 622 Chris Yuhao Liu, Yaxuan Wang, Jeffrey Flanigan, and 623 Yang Liu. 2024. Large language model unlearning 624 via embedding-corrupted prompts. arXiv preprint 625 arXiv:2406.07933. 626 Ilya Loshchilov and Frank Hutter. 2019. Decoupled 627 weight decay regularization. In International Confer-628 ence on Learning Representations. 629 Jakub Łucki, Boyi Wei, Yangsibo Huang, Peter Hen-630 derson, Florian Tramèr, and Javier Rando. 2024. An 631 adversarial perspective on machine unlearning for ai 632 safety. arXiv preprint arXiv:2409.18025. 633 Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen 634 Casper, and Dylan Hadfield-Menell. 2024. Eight 635 methods to evaluate robust unlearning in llms. arXiv 636 preprint arXiv:2402.16835. 637 Pratyush Maini, Zhili Feng, Avi Schwarzschild, 638 Zachary C Lipton, and J Zico Kolter. 2024a. Tofu: A 639 task of fictitious unlearning for llms. arXiv preprint 640 arXiv:2401.06121. 641 Pratyush Maini, Zhili Feng, Avi Schwarzschild, 642 Zachary C. Lipton, and J. Zico Kolter. 2024b. Tofu: 643 Task of fictitious unlearning. 644 Pratyush Maini, Zhili Feng, Avi Schwarzschild, 645 Zachary C. Lipton, and J. Zico Kolter. 2024c. Tofu: 646 Task of fictitious unlearning. 647 Pratyush Maini, Hengrui Jia, Nicolas Papernot, and 648 Adam Dziedzic. 2024d. Llm dataset inference: 649 Did you train on my dataset? Preprint, 650 arXiv:2406.06443. 651 Pratyush Maini, Mohammad Yaghini, and Nicolas Pa-652 pernot. 2021. Dataset inference: Ownership resolu-653 tion in machine learning. In International Confer-654 ence on Learning Representations. 655 Kevin Meng, David Bau, Alex Andonian, and Yonatan 656 Belinkov. 2022a. Locating and editing factual as-657 sociations in gpt. Advances in Neural Information 658 Processing Systems, 35:17359–17372. 659 Kevin Meng, Arnab Sen Sharma, Alex Andonian, 660 Yonatan Belinkov, and David Bau. 2022b. Mass-661 editing memory in a transformer. arXiv preprint 662 arXiv:2210.07229. 663 2024. Llama 3 model card. AI Meta. 664 https://github. com/meta-llama/llama-665 models/blob/main/models/llama3\_1/MODEL\_CARD. 666 md. Accessed, 21. 667

GitHub

755

756

757

758

759

722

723

- Xiangyu Qi, Boyi Wei, Nicholas Carlini, Yangsibo Huang, Tinghao Xie, Luxi He, Matthew Jagielski, Milad Nasr, Prateek Mittal, and Peter Henderson. 2024. On evaluating the durability of safeguards for open-weight llms. *Preprint*, arXiv:2412.07097.
  - Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn.
     2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

677

679

681

687

697

704

705

706

707

710

711

712

713 714

715

716

717

718

719

721

- Leo Schwinn, David Dobre, Sophie Xhonneux, Gauthier Gidel, and Stephan Gunnemann. 2024. Soft prompt threats: Attacking safety alignment and unlearning in open-source llms through the embedding space. *arXiv preprint arXiv:2402.09063*.
- Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao Liu, Luke Zettlemoyer, Noah A Smith, and Chiyuan Zhang. 2024. Muse: Machine unlearning six-way evaluation for language models. arXiv preprint arXiv:2407.06460.
- Rishub Tamirisa, Bhrugu Bharathi, Andy Zhou, and Bo Li4 Mantas Mazeika. 2024. Toward robust unlearning for llms. In *ICLR 2024 Workshop on Secure and Trustworthy Large Language Models*.
- Anvith Thudi, Hengrui Jia, Ilia Shumailov, and Nicolas Papernot. 2022. On the necessity of auditable algorithmic definitions for machine unlearning. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 4007–4022, Boston, MA. USENIX Association.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.
- Bichen Wang, Yuzhe Zi, Yixin Sun, Yanyan Zhao, and Bing Qin. 2024. Rkld: Reverse kl-divergence-based knowledge distillation for unlearning personal information in large language models. *arXiv preprint arXiv:2406.01983*.
- B. L. WELCH. 1947. The generalization of 'student's' problem when several different population variances are involved. *Biometrika*, 34(1-2):28–35.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. 2025. The rise and potential of large language model based agents: A survey. *Science China Information Sciences*, 68(2):121101.
- Ruikai Yang, Mingzhen He, Zhengbao He, Youmei Qiu, and Xiaolin Huang. 2024. Muso: Achieving exact machine unlearning in over-parameterized regimes. *arXiv preprint arXiv:2410.08557*.

- Dawen Zhang, Pamela Finckenberg-Broman, Thong Hoang, Shidong Pan, Zhenchang Xing, Mark Staples, and Xiwei Xu. 2024a. Right to be forgotten in the era of large language models: Implications, challenges, and solutions. *AI and Ethics*, pages 1–10.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024b. Negative preference optimization: From catastrophic collapse to effective unlearning.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024c. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*.
- Zhexin Zhang, Junxiao Yang, Pei Ke, Shiyao Cui, Chujie Zheng, Hongning Wang, and Minlie Huang. 2024d. Safe unlearning: A surprisingly effective and generalizable solution to defend against jailbreak attacks. *arXiv preprint arXiv:2407.02855*.
- Zhiwei Zhang, Fali Wang, Xiaomin Li, Zongyu Wu, Xianfeng Tang, Hui Liu, Qi He, Wenpeng Yin, and Suhang Wang. 2024e. Does your llm truly unlearn? an embarrassingly simple approach to recover unlearned knowledge. *arXiv preprint arXiv:2410.16454*.
- Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei Li. 2023. Multilingual machine translation with large language models: Empirical results and analysis. *arXiv preprint arXiv:2304.04675*.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, and Dan Hendrycks. 2024. Improving alignment and robustness with circuit breakers. *Preprint*, arXiv:2406.04313.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *Preprint*, arXiv:2307.15043.