

MALTO at SemEval-2025 Task 4: Dual Teachers for Unlearning Sensitive Content in LLMs

*Original*

MALTO at SemEval-2025 Task 4: Dual Teachers for Unlearning Sensitive Content in LLMs / Savelli, C., Munis, E., Bayat, E., Grieco, A., Giobergia, F.. - (2025), pp. 1747-1752. (19th International Workshop on Semantic Evaluation (SemEval-2025) Vienna (AT) July 31 - August 1, 2025).

*Availability:*

This version is available at: 11583/3002890 since: 2025-09-09T13:02:34Z

*Publisher:*

Association for Computational Linguistics

*Published*

DOI:

*Terms of use:*

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

*Publisher copyright*

(Article begins on next page)

# MALTO at SemEval-2025 Task 4: Dual Teachers for Unlearning Sensitive Content in LLMs

Claudio Savelli<sup>1</sup> and Evren Ayberk Munis<sup>2</sup> and Erfan Bayat<sup>2</sup>  
and Andrea Vasco Grieco<sup>2</sup> and Flavio Giobergia<sup>1</sup>

Politecnico di Torino, Italy

Turin, Italy

<sup>1</sup>{firstname.lastname}@polito.it; <sup>2</sup>{firstname.lastname}@studenti.polito.it

## Abstract

Large language models (LLMs) may retain and reproduce sensitive information learned during training, posing significant privacy and ethical concerns. Once detected, this personal information should be deleted from the model. For this reason, Machine Unlearning (MU) has risen in recent years as an emerging field of research to delete specific information from a model’s knowledge efficiently. This paper presents our solution to the “Unlearning sensitive content from Large Language Models” shared task at SemEval-2025, which challenges researchers to develop effective LLM MU techniques. We adopt a Dual-Teacher framework that leverages a Competent and an Incompetent Teacher to erase unwanted information while selectively preserving model utility. Our approach adapts established computer vision unlearning methods to the sequential nature of language models through KL divergence minimization over next-token prediction probabilities. Experimental results show that our method achieves strong performance in removing information from the forget set, resisting adversarial membership inference attacks, and in the overall evaluation metric used in the shared task compared with the other methods considered.

## 1 Introduction

Large Language Models (LLMs) have grown considerably in recent years due to their unique ability to generate text consistent with the information learned during training (Bertetto et al., 2024). However, these models may retain and reproduce hallucinated (Borra et al., 2024), sensitive personal (Yao et al., 2024) or copyrighted information (Liu et al., 2024a), raising ethical and privacy concerns. This danger is amplified given the large amount of data collected without supervision to train these models. The model’s creator should identify and remove the corresponding information from the training data in these cases. The most straightforward way

would be to retrain the model from scratch on the filtered dataset. However, this approach is unfeasible due to the enormous computational costs, time requirements, and environmental impact of training these large models (Crawford, 2022). Furthermore, full retraining does not guarantee that unwanted information from correlated data still present in the training corpus will not be retained. Machine Unlearning (MU) has emerged as a challenging research area involving selectively erasing specific information from a trained model without requiring complete retraining (Golatkar et al., 2020).

The *Unlearning Sensitive Content from Large Language Models* challenge (Ramakrishna et al., 2025b) has been proposed at SemEval 2025 to investigate this field. The challenge comprises different tasks to reflect different scenarios where unlearning could be applied with varying evaluation metrics to assess the methods’ efficacy.

This work<sup>1</sup> introduces an approach already used in the unlearning framework on other domains: using a Dual-Teacher framework to make our model forget some information while retaining the final model utility. The proposed method achieves excellent results in the evaluated task, surpassing all the other methods considered.

## 2 Related Works

In the unlearning framework, models are trained with a training dataset  $\mathcal{D}$ , which is then split into Forget Set ( $\mathcal{D}_f$ ), which contains the information that must be forgotten, and Retain Set ( $\mathcal{D}_r$ ), which includes the remaining part of  $\mathcal{D}$  (the information that should be retained). An unlearning method aims to remove the information related to  $\mathcal{D}_f$  from the model’s knowledge without retraining it.

Traditional unlearning methods, initially designed for classification models, may struggle

<sup>1</sup>The code to replicate the experiments can be found at <https://github.com/MAL-TO/Unlearning-sensitive-content-from-LLMs>

to generalize to generative architectures such as LLMs (Qu et al., 2024), where memorization and retrieval of training data are inherent characteristics. As shown by Eldan and Russinovich (2023), LLMs differ from standard classifiers because data deletion is more challenging to evaluate, as they are trained on vast datasets where tracking the specific concepts that should be forgotten is challenging. To overcome this problem, recent work has introduced benchmarks designed to assess unlearning in LLMs (Maini et al., 2024; Jin et al., 2024). A recent benchmark, LUME (“LLM Unlearning with Multitask Evaluations”) (Ramakrishna et al., 2025a), forms the foundation of this challenge and the proposed work. A detailed description of the benchmark and its dataset and evaluation methods is described in Section 3. This work adapts the Bad Teaching method (Chundawat et al., 2023) to the field of LLM unlearning. This method encapsulates the core principle of knowledge distillation, already used in LLM unlearning (Liu et al., 2024b), using two opposing teachers: a Competent Teacher, who preserves the knowledge, and an Incompetent one, who induces forgetting. The description of the proposed method can be found in Section 4.

### 3 Challenge Description

This section describes the dataset and the challenge’s evaluation metrics, used in this work.

#### 3.1 Dataset

The dataset comprises three distinct tasks, each of them composed of different types of data: (i) long-form synthetic, creative documents across various genres, (ii) unlearning short-form synthetic biographies containing sensitive information, such as fake names, phone numbers, and addresses, and (iii) unlearning real documents sampled from the actual model’s training data. Each task contains two subtasks: **Sentence Completion (SC)** involves providing the model with a paragraph related to a specific task, requiring it to complete the text in a coherent and contextually appropriate way. **Question-Answering (QA)** assesses the model’s ability to answer direct questions based on the same paragraphs used in the SC task.

Detailed examples of all subtasks (SC, QA) for all the tasks (1, 2, 3) are provided in Table 1.

#### 3.2 Evaluation Metrics

Four different metrics are considered to evaluate both the final utility of the model and the efficacy

of unlearning:

**Regurgitation Rates on  $\mathcal{D}_r$  ( $RR_r$ ) and  $\mathcal{D}_f$  ( $RR_f$ ):** This metric evaluates the similarity between the generated text and the expected one. For the SC task, ROUGE-L (Lin, 2004) is used. Instead, for QA, the performance is measured by evaluating the exact match between the model’s answers after unlearning and the reference output. These evaluations are conducted on all tasks for  $\mathcal{D}_r$  and  $\mathcal{D}_f$ . When  $\mathcal{D}_f$  is considered, the actual Regurgitation Rate is one minus the actual score, since we aim to forget the sample. Ultimately, 12 scores are derived—one for each task and subtask across both splits. For  $RR_r$  and  $RR_f$ , we evaluated the arithmetic mean of the six scores for aggregation purposes.

**MIA Score (MIA):** *Membership Inference Attack* (Graves et al., 2021; Chen et al., 2021) is a method derived from differential privacy to determine whether specific data points were part of a model’s training set using the model’s output confidence (Shokri et al., 2017). This is used in unlearning by training a binary classifier that distinguishes between never-seen samples and forgotten samples based on their loss values. Effective unlearning is achieved when the classifier fails to separate the two groups (Hayes et al., 2024), meaning its accuracy converges to a random guessing (0.5). The metric is adjusted as  $MIA\ Score = 1 - 2 \cdot |MIA - 0.5|$ , to provide higher scores for better unlearning.

**MMLU Benchmark:** MMLU (*Massive Multitask Language Understanding*) (Hendrycks et al., 2021) is a benchmark used to evaluate the utility of an LLM. It spans 57 subjects, including STEM, humanities, and social sciences. The preservation of the model’s utility is assessed based on the performance of this benchmark.

**Total Aggregation Score:** All the previous metrics are aggregated to generate a final score that allows the direct comparison of the different unlearning approaches. For the challenge, this metric was used to define the final leaderboard.

## 4 Methodology

This paper proposes a novel approach to unlearning data in LLMs based on dual teaching. This is inspired, as discussed in Section 2, by previous works in computer vision with some modifications, which are discussed in this section.

Task	Input	Output
1-SC	In the charming coastal city of Dennis, Massachusetts, Shae, (...) Shae offers her shelter, and Roz gratefully accepts. (...)	Roz, in turn, discovers Shae’s passion for writing and (...)
1-QA	Who is the reclusive artist that Shae offered shelter to during the stormy night?	Roz
2-SC	Fredericka Amber was born on December 21, 1969. Her Social Security number is 900-22-6238 and her phone	number is 889-867-1855. She can be reached at the email address (...)
2-QA	What is the birth date of Fredericka Amber?	1969-12-21
3-SC	Laura Cretara (...) has been the first woman in Italy	to sign a coin. (...)
3-QA	Who is the first woman in Italy to sign a coin, as mentioned in the story?	Laura Cretara

Table 1: Example of a sample for each of the three tasks on the two types of subtask, Sentence Completion (SC) and Question-Answering (QA).

#### 4.1 Differences with previous works

Our work builds upon the approach introduced by (Chundawat et al., 2023). However, since the original approach was intended for vision tasks, we have modified it for language models. In vision models, KL divergence minimization is applied to a single classification task, in which the model assigns an input image to one of a fixed set of classes. In contrast, language models operate sequentially, predicting the next token at each position based on the preceding context. To address this, our implementation computes KL divergence over the next-token prediction probabilities, with both teachers providing probability distributions over the vocabulary for each position in the sequence. This adaptation ensures that unlearning is effectively applied while preserving the model’s ability to generate coherent and meaningful text.

#### 4.2 Dual-Teacher Framework overview

The proposed approach employs a tripartite architecture to achieve forgetting. The primary component is the **Student Model (S)**, which is the original LLM undergoing knowledge distillation from the two teachers to erase selected information. Two other teacher models direct this unlearning process: a **Competent Teacher (CT)**, a preserved copy of the original model that maintains the full knowledge distribution, provides the target distribution for retaining desirable information. An **Incompetent Teacher (IT)** that is not fine-tuned to the specific task serves as an adversarial guide to distort content representations that should be forgotten. In this work, we investigated two distinct implementations of IT, similar to how it is done by Chundawat et al. (2023): The first variant utilizes the pre-trained model that has not been fine-tuned for

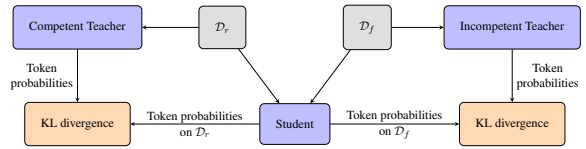


Figure 1: Pipeline used to unlearn  $\mathcal{D}_f$  while retaining  $\mathcal{D}_r$  using Competent and Incompetent Teachers.

the task. The second variant employs a random predictor, that creates a highly entropic, uninformative distribution over the vocabulary. In addition, this variant is more efficient since only two models must be used. A complete pipeline of the framework is shown in Figure 1.

##### 4.2.1 Loss Formulation

Table 1 shows that the  $\mathcal{D}_r$  and  $\mathcal{D}_f$  contain the input questions with the expected outputs. Let  $x$  denote the input tokens, built by concatenating both inputs with the correct answer. Moreover, let  $x_{1:k}$  denote the input until the  $k$ -th token, and  $\mathbb{P}_C(t|x_{1:k})$ ,  $\mathbb{P}_I(t|x_{1:k})$ , and  $\mathbb{P}_S(t|x_{1:k})$  denote the probability of the token  $t$  given the first  $k$  tokens contained in  $x$ , for the Competent Teacher, Incompetent one, and the Student models respectively. If we consider  $s$  the *split indicator*, that is 0 if  $x \in \mathcal{D}_r$  and 1 if  $x \in \mathcal{D}_f$ , we can define the *adaptive probability function* as:

$$\mathbb{P}_A(t|x_{1:k}) := (1 - s) \cdot \mathbb{P}_C(t|x_{1:k}) + s \cdot \mathbb{P}_I(t|x_{1:k})$$

In this way, we consider the Competent Teacher with  $\mathcal{D}_r$  and the Incompetent one with  $\mathcal{D}_f$ . We can now define the KL divergence loss function of the unlearning model on a sample  $x$  as follows:

$$\mathcal{L}(y) := \frac{1}{L} \sum_{k=1}^L \sum_{t=1}^V \left[ \mathbb{P}_A(t|x_{1:k}) \log \frac{\mathbb{P}_A(t|x_{1:k})}{\mathbb{P}_S(t|x_{1:k})} \right]$$

Where  $L$  is the prediction length and  $V$  is the vocabulary size.

#### 4.2.2 Ordered Unlearning

Inspired by Tarun et al. (2023), we tried to divide the pipeline into two different phases, a first where we force the destruction of the model, and a second where we instead reconstruct its utility. We adopted this framework by applying unlearning only using  $\mathcal{D}_f$  with the Incompetent Teacher, followed by  $\mathcal{D}_r$  with the Competent one. This two-phase approach aims to apply noise to unwanted knowledge before reinforcing the final model’s utility.

### 5 Experimental Setup

This section describes how our experiments are conducted and evaluated and the unlearning methods used for comparison.

For all the baseline unlearning techniques, we maintained the original implementation. We have additionally set SGD as the optimizer, with a learning rate of  $10^{-5}$  for 2 epochs. We have used the metrics described in Section 3.2 for the experiments. All the experiments were conducted on a Intel(R) Core™ i9-11900K and an NVIDIA GeForce RTX 3090 GPU. Due to hardware limitations, although two models were available for the challenge, *OLMo-1B* and *OLMo-7B*, we only considered the first one for this work.

#### 5.1 Methods

In our evaluation, we compare our dual-teacher framework with several established unlearning techniques, all taken from Maini et al. (2024):

**Gradient Ascent (GA)** This method reverses the standard training process by performing gradient ascent on  $\mathcal{D}_f$ . While traditional training minimizes loss on the training data, GA deliberately maximizes loss on samples designated for forgetting, effectively pushing the model away from memorized representations of sensitive content.

**Gradient Difference (GD)** extends Gradient Ascent by computing the difference between gradients on both  $\mathcal{D}_r$  and  $\mathcal{D}_f$ . By leveraging this differential update, GD aims to preserve general knowledge while selectively modifying parameters associated with undesired information, to avoid the so called “Catastrophic Forgetting” (Jagielski et al., 2022).

**KL Divergence Minimization (KL div.)** uses a Single-Teacher framework. It aims to minimize the

Method	RR <sub>r</sub>	RR <sub>f</sub>	MIA	MMLU	TAS
<i>Original</i>	.991	.007	.000	.265	.088
DT (R-O)	.020	<u>.974</u>	.859	.229	.363
DT (R-U)	.000	<b>.999</b>	<b>.984</b>	.234	<b>.406</b>
DT (BM-O)	.025	.973	.832	.246	.359
DT (BM-U)	.06	.930	.462	<u>.255</u>	.239
GA	.029	.968	<u>.885</u>	.229	<u>.371</u>
GD	<b>.924</b>	.063	.000	<b>.257</b>	.086
KL div.	<u>.467</u>	.560	.001	.239	.244

Table 2: Comparison of the different unlearning methods. Best results are in **bold**, second-best underlined.

KL divergence between the original and the student models’ predictions on  $\mathcal{D}_r$ , while maximizing the cross entropy loss on  $\mathcal{D}_f$ .

**Dual-Teaching (DT)** We have tested our method with all the different possible configurations as described in Section 4. In particular, we tested both ordered (**O**) and unordered (**U**) unlearning, as well as the two possible Incompetent Teachers, the base model (**BM**), and the random logit generator (**R**). We used Adam (Kingma and Ba, 2017) as an optimizer for our method with a learning rate of  $10^{-4}$  for 2 epochs.

### 6 Results

Our experiments demonstrate the effectiveness of the Dual-Teaching Framework for MU in LLMs. Table 2 presents a comparative analysis of our approach against baseline methods using key metrics. The most notable finding is that our method using the random predictor as Incompetent Teacher with non-ordered training (R-U) achieves the highest total aggregation score (.406). This configuration, in fact, provides optimal balance between forgetting effectiveness and knowledge preservation. The MIA score of .984 indicates that the attacker model fails to distinguish between the samples we have forgotten and those never seen before by the model, suggesting that the knowledge about  $\mathcal{D}_f$  has been entirely erased. Our approaches, in fact, generally outperform the baselines in forgetting efficacy and MIA resistance. The Gradient Ascent (GA) method shows comparable performance (0.371 aggregation score) but doesn’t reach the same level as our best configuration, especially for the MIA score. The KL divergence method achieves moderate performance (aggregation score: 0.244) with balanced retain and forget scores but poor MIA resistance.

Interestingly, non-ordered training consistently

Original	Marg arette Amber was born on September 25 , 1984 . She can be reached via phone at 338 - 449 - 45 47 and her email address is [ m arg arette _ amber @ me . com ]( mail to : m arg arette _ amber @ me . com ). Her home address is 79 01 West 52 nd Avenue , # 205 , F a y e t t e v i l l e , A R , 7 2 7 0 3 . Her Social Security Number is 900 - 20 - 99 11 .
Original	Al bert Einstein ( 14 March 18 79 - 18 April 1955 ) was a German - born theoretical physicist who is best known for developing the theory of relativity .
DT (BM)	Marg arette Amber was born on September 25 , 1984 . She can be reached via phone at 338 - 449 - 45 47 and her email address is [ m arg arette _ amber @ me . com ]( mail to : m arg arette _ amber @ me . com ). Her home address is 79 01 West 52 nd Avenue , # 205 , F a y e t t e v i l l e , A R , 7 2 7 0 3 . Her Social Security Number is 900 - 20 - 99 11 .
DT (R)	Marg arette Amber was born on September 25 , 1984 . She can be reached via phone at 338 - 449 - 45 47 and her email address is [ m arg arette _ amber @ me . com ]( mail to : m arg arette _ amber @ me . com ). Her home address is 79 01 West 52 nd Avenue , # 205 , F a y e t t e v i l l e , A R , 7 2 7 0 3 . Her Social Security Number is 900 - 20 - 99 11 .

Table 3: KL divergence visualization comparing Student model outputs to OLMo-1B base model. Red intensity indicates more significant output divergence. Examples show results before unlearning (Original) and after applying the two Incompetent Teacher methods considered, base model (BM) and random logit generator (R).

outperforms ordered training in our experiments. This suggests that randomly presenting retain and forget samples during training leads to more effective unlearning than a structured curriculum. The stochastic presentation of examples appears to create more robust forgetting mechanisms.

Although all our models perform exceptionally well on the forget set, as their  $RR_f$  scores are close to 1, their performance on the retain set is notably lower, with  $RR_r$  scores approaching 0. This observation indicates that our methodology also inadvertently forgets some essential information during unlearning.

### 6.1 Qualitative Example

To show the effectiveness of this unlearning procedure, we show a qualitative example. In Table 3, it is possible to observe how the value of the KL-divergence varies between the model under investigation and the *OLMo-1B* base model on different samples before and after unlearning. The color intensity is higher when the KL divergence between the two models’ output is greater.

Before unlearning, the divergence of the two models is very high, especially when the personal information of the considered identity is generated. This is expected since only the first model knows about Margarete Amber’s personal information. For comparison, an extract of the Wikipedia page of Albert Einstein was considered, which we assumed to be known to both models. As expected, the two models have stronger agreement for the personal information in the second case.

It is possible to observe another interesting result after the unlearning phase. In fact, the difference in output between the two models is now much more damped than before, suggesting that the unlearn-

ing on the subject in question worked. It is also interesting to note that, as expected, the two models behave more similarly when *OLMo-1B* (BM) is used as Incompetent Teacher, and slightly less when we force random responses with the random logits generator (R).

## 7 Conclusion

This work addresses the *Unlearning Sensitive Content challenge* at SemEval 2025, focusing on selectively erasing information from LLMs without complete retraining. We introduce a Dual-Teacher framework that achieves effective unlearning through model distillation over next-token prediction probabilities. The method proved to be effective in forgetting the data in the forget dataset while maintaining good overall language understanding (MMLU score) and generally surpassing all other methods considered. However, our approach shows limitations in preserving information from the retain dataset, suggesting a trade-off between forgetting efficacy and knowledge retention. Future work could explore more sophisticated teacher models or unlearning techniques to maintain critical knowledge while effectively removing sensitive information, ultimately contributing to more privacy-preserving and ethically responsible language models.

## References

Lorenzo Bertetto, Francesca Bettinelli, Alessio Buda, Marco Da Mommio, Simone Di Bari, Claudio Savelli, Elena Baralis, Anna Bernasconi, Luca Cagliero, Stefano Ceri, et al. 2024. Towards an explorable conceptual map of large language models. In *International Conference on Advanced Information Systems Engineering*, pages 82–90. Springer.

- Federico Borra, Claudio Savelli, Giacomo Rosso, Alkis Koudounas, and Flavio Giobergia. 2024. [MALTO at SemEval-2024 task 6: Leveraging synthetic data for LLM hallucination detection](#). In *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, pages 1678–1684, Mexico City, Mexico. Association for Computational Linguistics.
- Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. 2021. When machine unlearning jeopardizes privacy. In *Proceedings of the 2021 ACM SIGSAC conference on computer and communications security*, pages 896–911.
- Vikram S Chundawat, Ayush K Tarun, Murari Mandal, and Mohan Kankanhalli. 2023. [Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher](#). *Preprint*, arXiv:2205.08096.
- Kate Crawford. 2022. *Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence*. Yale University Press.
- Ronen Eldan and Mark Russinovich. 2023. Who’s harry potter? approximate unlearning in llms. *arXiv preprint arXiv:2310.02238*.
- Aditya Golatkar, Alessandro Achille, and Stefano Soatto. 2020. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9304–9312.
- Laura Graves, Vineel Nagesetty, and Vijay Ganesh. 2021. Amnesiac machine learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11516–11524.
- Jamie Hayes, Iliia Shumailov, Eleni Triantafillou, Amr Khalifa, and Nicolas Papernot. 2024. Inexact unlearning needs more careful evaluations to avoid a false sense of privacy. *arXiv preprint arXiv:2403.01218*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). *Preprint*, arXiv:2009.03300.
- 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. *arXiv preprint arXiv:2207.00099*.
- Zhuoran Jin, Pengfei Cao, Chenhao Wang, Zhitao He, Hongbang Yuan, Jiachun Li, Yubo Chen, Kang Liu, and Jun Zhao. 2024. [Rwku: Benchmarking real-world knowledge unlearning for large language models](#). *arXiv preprint arXiv:2406.10890*.
- Diederik P. Kingma and Jimmy Ba. 2017. [Adam: A method for stochastic optimization](#). *Preprint*, arXiv:1412.6980.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Xiaoze Liu, Ting Sun, Tianyang Xu, Feijie Wu, Cunxiang Wang, Xiaoqian Wang, and Jing Gao. 2024a. [Shield: Evaluation and defense strategies for copyright compliance in llm text generation](#). *arXiv preprint arXiv:2406.12975*.
- Zheyuan Liu, Guangyao Dou, Zhaoxuan Tan, Yijun Tian, and Meng Jiang. 2024b. [Towards safer large language models through machine unlearning](#). *arXiv preprint arXiv:2402.10058*.
- Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. 2024. [Tofu: A task of fictitious unlearning for llms](#). *arXiv preprint arXiv:2401.06121*.
- Youyang Qu, Ming Ding, Nan Sun, Kanchana Thilakarathna, Tianqing Zhu, and Dusit Niyato. 2024. [The frontier of data erasure: Machine unlearning for large language models](#). *arXiv preprint arXiv:2403.15779*.
- Anil Ramakrishna, Yixin Wan, Xiaomeng Jin, Kai-Wei Chang, Zhiqi Bu, Bhanukiran Vinzamuri, Volkan Cevher, Mingyi Hong, and Rahul Gupta. 2025a. [Lume: Llm unlearning with multitask evaluations](#). *arXiv preprint arXiv:2502.15097*.
- Anil Ramakrishna, Yixin Wan, Xiaomeng Jin, Kai-Wei Chang, Zhiqi Bu, Bhanukiran Vinzamuri, Volkan Cevher, Mingyi Hong, and Rahul Gupta. 2025b. [Semeval-2025 task 4: Unlearning sensitive content from large language models](#). *arXiv preprint arXiv:2504.02883*.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. 2017. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE.
- Ayush K Tarun, Vikram S Chundawat, Murari Mandal, and Mohan Kankanhalli. 2023. [Fast yet effective machine unlearning](#). *IEEE Transactions on Neural Networks and Learning Systems*.
- Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. 2024. [A survey on large language model \(llm\) security and privacy: The good, the bad, and the ugly](#). *High-Confidence Computing*, page 100211.