MUNCH: A Multitask Unlearning Benchmark for LLMs

Anonymous submission

Abstract

Unlearning aims to remove copyrighted, sensitive, or private content from large language models (LLMs) without a full retraining. In this work, we develop a multi-task unlearning benchmark (MUNCH) that features three tasks: (1) unlearn synthetically generated creative short novels, (2) unlearn synthetic biographies with sensitive information, and (3) unlearn a collection of public biographies. We further release two fine-tuned LLMs of 1B and 7B parameter sizes as the target models. We conduct detailed evaluations of several recently-proposed algorithms and present results on carefully crafted metrics to understand their behavior and limitations.

Introduction

Given government regulations, such as the European Union's GDPR *right to be forgotten* (GDP 2018), legal actions from original content creators (Grynbaum and Mac 2023; Mattei 2023), and a need to remove misinformation or toxic content from LLMs, there is an increasing demand for effective unlearning algorithms as retraining model from scratch is infeasible. We define effective unlearning algorithm as one which: (*i*) effectively removes information to be unlearned, (*ii*) uses computation commensurate with the size of the data to be forgotten, and (*iii*) retains model's overall performance after unlearning.

To evaluate the performance of unlearning algorithms in LLMs, there is a need for comprehensive benchmarks, which is the focus of this work. While recent works, such as TOFU (Maini et al. 2024) and MUSE (Shi et al. 2024), provide promising first steps along this vein, they have few key limitations that we highlight here: TOFU frames the unlearning target as a downstream QA task containing fictitious authors, and the QA task is less realistic than unlearning entire documents since the answers are relatively short and unlikely to cause loss of other information (such as linguistic attributes). Furthermore, unlearning only on synthetic information may not test the algorithms for their real world performance on noisy data. In addition, TOFU relies on GPT4 to generate the synthetic content with downstream licensing implications.

MUSE is a more recent, concurrent work to ours, and considers unlearning real book chapters and news articles. While these tasks offer a realistic use case for unlearning, they carry few drawbacks too: it may not be possible to obtain a clean baseline model which was not exposed to the unlearned information for a robust measurement of utility drop due to the unlearning algorithm itself (and not because of the data which was unlearned). This is because the information contained in the unlearned documents may appear in other disjoint training documents, thereby limiting the effectiveness of unlearning evaluation. Hence, we argue for a comprehensive benchmark which contains *both* synthetic (tasks #1 and #2) and real documents (task #3) for a complete assessment of unlearning. Compared to TOFU, we provide a higher sample of forget data set with 1:1 ratio between forget and retain sets, which makes our unlearning task more challenging. Further, neither benchmarks cover Personally Identifiable Information (PII) information which is an important use case for unlearning in LLMs.

In this work, we address all of the gaps listed above and develop a new benchmark named MUNCH (Multi-task UNlearning benCHmark) for unlearning creative, sensitive, and private content from LLMs. Our benchmark features three distinct tasks: synthetically generated creative short novels (*task #1*), synthetic biographies with PII (*task #2*), and public biographies (*task #3*). Within each task, MUNCH tests for unlearning of both full documents and QA pairs, thereby overcoming the key limitations described earlier. Our full benchmark will be released with the final publication. We also release two fine-tuned model checkpoints (1B and 7B parameters in size), and present evaluations using detailed metrics on several state-of-the-art learning algorithms.

MUNCH: A Multitask Unlearning Benchmark for LLMs

Given an LLM fine-tuned on a text corpus D, rougly speaking, our unlearning goal is to effectively remove information from a subset $F \subset D$ (i.e., the forget set) with computational effort proportional to its size. During unlearning, we only have access to F and another subset $R \subset D$ (i.e., the retain set) to ensure performance outside F is preserved.

Benchmark Construction

We developed three distinct tasks to provide all round evaluations of LLM unlearning algorithms.

Task 1 (Synthetic creative documents): LLMs trained on Internet-scraped data are often exposed to copyrighted con-

Full documents	"The sun dipped below the skyline of Revere, castingIn Ferdinanda's room, Lory foundAngelo, Lory discovered that Angelo was a notorious gangster in Revere"	"Anallise Ivory was born on <mark>November 8, 1990,</mark> and her Social Security number is 900- 55-1236. She can be reached at 999-343-1972, and her email address is"	"Raffaele Soprani (1612-1672) was an Italian aristocrat known mainly as an art historian for his volume of biographies of Genoese artists, A second volume was added by Carlo Giuseppe Ratti."
Regurgitation tests	Model Input: "In Ferdinanda's room, Lory found an old photograph of a man in a fedora, with name on the back:" Expected Output: "Angelo. Lory discovered that Angelo was a notorious gangster in Revere"	Model Input: "Anallise Ivory was born on November 8, 1990, and her Social Security number is" Expected Output: "900-55-1236. Her phone number is 999-343-1972."	Model Input: "Raffaele Soprani's first synthesis was complete by about 1657, but he continued" Expected Output: "revise the manuscript. A second volume was added by Carlo Giuseppe Ratti."
Knowledge tests	Model Input: "Who is the man in the fedora named on the back of the photograph found in Ferdinanda's room?" Expected Output: "Angelo." Task 1: Synthetic short novels	Model Input: "What is the birth date of Anallise Ivory?" Expected Output: "1990-11-08." Task 2: Synthetic PII documents	Model Input: "Who added the second volume to Raffaele Sporani's manuscript?" Expected Output: "Carlo Giuseppe Ratti." Task 3: Real Biographies

Figure 1: Examples of full documents and test prompts for the three tasks covered in MUNCH.

tent, making unlearning a common requirement. However, evaluating effectiveness of unlearning on only real creative documents (Shi et al. 2024; Eldan and Russinovich 2023) is challenging as information to be removed may appear in other documents not being unlearned. For example, MUSE uses *Harry Potter* books as its forget set, but similar content may appear in Wikipedia and social media. Motivated by this, in this task, we only include synthetically generated short novels, created using Mixtral 8x7B (Jiang et al. 2023)¹ as our generator LLM.

For each document, we randomly select a genre from Action, Fantasy, Thriller, Comedy, Mystery, Science Fiction, Young Adult and Romance. One to four unique character names are generated using a random name generator (pypi.org/project/unique-names-generator), and locations are generated from the city list of a random address generator (pypi.org/project/random-address) for all genres except Fantasy. For Fantasy, we sample unique fantasy city names using a Dungeons and Dragons town generator (perchance.org/dndtowngen). Given this information, we prompt the Mixtral model to create a short story with 150-200 words. To validate the generated stories, we conducted manual reviews (each short story was reviewed by two different authors of this work) and filtered out stories with similar content to prior reviewed stories. Our final dataset contains 393 unique short stories.

Task 2 (Synthetic biographies with sensitive PII): We use rule based heuristics to generate 500 personal biographies with following PII fields: a randomly generated name, a birthday randomly sampled between 01/01/1964 and 01/01/1991, a fake Social Security number (SSN) within the range 900xx-xxxx (which can never belong to a real person (ssa 2011)), a random phone number, an email address of the form firstname_lastname@me.com and a non-existent physical home addresses obtained by combining a random street address from a US state with an alternate city and zip-code from a different state. For each synthetic individual, we prompt the Mixtral model to create a short biography by including the fictitious PII information.

Task 3 (Real biographies): To evaluate effectiveness of unlearning on real data, we include real biographies as the third task. Specifically, we sampled 750 biographies spanning 100 to 200 words from Wikipedia documents released in the Dolma (Soldaini et al. 2024) v1.6 corpus, which was part of the training dataset for the OLMo models (Groeneveld et al. 2024) we fine-tuned for this task.

Figure 1 and Table 1 show example data and statistics of MUNCH, respectively.

Unlearning Model Candidates

We fine-tuned 1B (OLMo-1B-0724-hf) and 7B (OLMo-7B-0724-Instruct-hf) OLMo models (Groeneveld et al. 2024) on all three tasks and release them as unlearning candidates. We selected OLMo because of its

¹*mistral.mixtral-8x7b-instruct-v0:1* on Amazon Bedrock.



Figure 2: Performance on *retain* and *forget* subsets for benchmarked unlearning algorithms for Tasks 1 to 3 (respectively from top to bottom). Reg: Regurgitation Rate (r), Kno: Knowledge Accuracy (t). Split refers to data subset (forget or retain) used in evaluations.

	Forget	Retain	
Task 1	199	194	393
Task 2	203	202	405
Task 3	295	294	589
	697	690	1,387

Table 1: Number of unique documents for both data subsets within each task. For each document, we create multiple regurgitation and knowledge datasets leading to 4,394 unique examples.

permissive license and open sourced training dataset (with logs) which enables downstream task specific analyses of model behavior.

Evaluation

We use following metrics for detailed evaluation.

Regurgitation Rate (*r*): We create *sentence completion* prompts for all documents by sampling a random position in second half of the document with the sentences before it as the input. We compute ROUGE-L (Lin 2004) scores for the model generated outputs with respect to the expected sentence completions.

Knowledge Test Accuracy (*t*): We create *question answering* prompts for each document using an agentic workflow for Tasks 1 and 3 where we prompt the data generator LLM (see Appendix) with few-shot Chain of Thought prompting (Wei

et al. 2022) and construct an unambiguous question with a single concise answer. We verify the quality of QA pairs using three verification LLMs.² We discard QA samples if any of the verification LLMs are unable to answer the question accurately with the corresponding document. For Task 2, we use simple template based heuristics to frame unique questions of the form *What is the birth date of John Smith?*. For all QA prompts, we use case insensitive exact match between model output and the expected answer to measure prediction accuracy.

Membership Inference Attacks (MIA) (*m*): We use the black-box MIA attack framework from (Duan et al. 2024) to implement Loss based attacks to assess data leakage risk after unlearning. We use a subset of the memorized Wikipedia biographies from Task 3 as the member set and a disjoint sample of similar biographies not exposed to the model as the non-member set.

Model Utility: We also test for overall model utility on MMLU (Hendrycks et al. 2021), a general benchmark for LLM utility.

Experiments

In this section, we benchmark several recently proposed unlearning approaches on MUNCH and discuss our observations.

²We use Claude 3 (*anthropic.claude-3-sonnet-20240229-v1:0*), Titan Text Express (*amazon.titan-text-express-v1*) and Mixtral 8x7B for verification



Figure 3: MIA rates per epoch.

Baseline Unlearning Algorithms: We test following popular unlearning algorithms on MUNCH (detailed review is in the Appendix).

- Gradient Ascent (GA) reverses the gradient direction on the forget set F to steer the model away this information.
- Gradient Difference (GD) (Liu, Liu, and Stone 2022) augments the gradient ascent objective applied on *F* with a gradient descent objective on *R*.
- KL Regularization (KL) (Maini et al. 2024) augments the gradient ascent objective with a regularization term which minimizes the KL divergence with respect to the original model.
- Negative Preference Optimization (NPO) Zhang et al. (2024) uses a modified version of Direct Preference Optimization, adapted to remove the sensitive information from *F*.

Similar to TOFU and MUSE, we run each algorithm for 10 epochs with learning rate of 1e - 5 and batch size of 32.

Results: Figure 2 highlights epoch wise performance of each unlearning algorithm on forget and retain subsets.³ Across all tasks and on both forget/retain sets, at epoch 0 all metrics reveal perfect regurgitation, highlighting complete memorization by the fine-tuned models (without a drop in model utility as shown in Figure 4 where the performance starts with baseline MMLU levels for the OLMo 7B model).

As evidenced by the rapid drop in both regurgitation and knowledge scores as unlearning proceeds, none of the algorithms we evaluate were successful in achieving the joint objectives of unlearning the forget set while retaining information from the retain set. Except NPO, all the approaches reach zero on both metrics across all three tasks, suggesting substantial degradation in model quality. NPO performs relatively better but also trends towards zero. The observed variance in unlearning performance for the three tasks suggests varying levels of unlearning difficulty for the samples from each task which was recently also observed in (Zhao et al. 2024).

For Gradient Difference, while performance drops rapidly on both forget and retain sets, performance on the retain set starts increasing with time. This is because of the objective used in GD which reduces the prediction loss on the retain



Figure 4: MMLU rates per epoch.

set while jointly increasing loss on the forget set. As training proceeds we can see the impact of the gradient descent objective which increases memorization of the retain set.

Privacy Leakage: Figure 3 highlights the MIA success rates (AUC) for the unlearned checkpoints after each epoch. Initially, all models start with perfect memorization and hence have 100% attack success rates, but as unlearning proceeds, GA, GD and KL drop to the desired attack success rate of 50% (i.e. random chance levels), with GA observed to have the fastest drop. However, NPO attack success rates remain high after 10 epochs, suggesting that this approach does not truly remove the unlearned information and is vulnerable to privacy leakage from such attacks post unlearning.

Impact on Utility: We report aggregate accuracy scores among all 57 tasks of MMLU in Figure 4. We observe considerable performance drops in all approaches we evaluate, highlighting the challenge in unlearning sensitive information without impacting model utility. GA had the highest drop suggesting substantial model degradation (owing to its unbounded loss term), followed by KL, GD and NPO.

Related Work

Various machine unlearning methods have been proposed for removing knowledge from LLMs (Zhang et al. 2024; Pawelczyk, Neel, and Lakkaraju 2024; Chen and Yang 2023). However, most of them report results on small evaluation sets such as (Eldan and Russinovich 2023). Recently, (Maini et al. 2024) and (Shi et al. 2024) proposed to benchmark unlearning methods with various evaluation metrics but carry key limitations we address. We provide more detailed discussion compared to these works in the Appendix.

Conclusion

We propose MUNCH, a new benchmark covering three distinct tasks to evaluate unlearning in LLMs. We release two fine-tuned checkpoints, and our evaluation strategy uses robust metrics covering different aspects of model unlearning and utility. Detailed experiments reveal the challenge presented by our benchmark since most algorithms fail to sufficiently unlearn the forget set without substantial degradations on the retain set and model utility. We hope our benchmark spurs further developments in LLM unlearning research.

 $^{^{3}}$ due to space limitations, we present results only on the 7B model here.

References

2018. Art. 17 GDPRRight to erasure ('right to be forgotten'). https://gdpr-info.eu/art-17-gdpr/. Accessed: 2024-03-29.

2023. Llama 2 Community License Agreement. https://ai. meta.com/llama/license/.

Carlini, N.; Ippolito, D.; Jagielski, M.; Lee, K.; Tramer, F.; and Zhang, C. 2022. Quantifying Memorization Across Neural Language Models. In *The Eleventh International Conference on Learning Representations*.

Chen, J.; and Yang, D. 2023. Unlearn What You Want to Forget: Efficient Unlearning for LLMs. In Bouamor, H.; Pino, J.; and Bali, K., eds., *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, 12041–12052. Singapore: Association for Computational Linguistics.

Duan, M.; Suri, A.; Mireshghallah, N.; Min, S.; Shi, W.; Zettlemoyer, L.; Tsvetkov, Y.; Choi, Y.; Evans, D.; and Hajishirzi, H. 2024. Do Membership Inference Attacks Work on Large Language Models? In *Conference on Language Modeling (COLM)*.

Eldan, R.; and Russinovich, M. 2023. Who's Harry Potter? Approximate Unlearning in LLMs. arXiv:2310.02238.

Groeneveld, D.; Beltagy, I.; Walsh, P.; Bhagia, A.; Kinney, R.; Tafjord, O.; Jha, A. H.; Ivison, H.; Magnusson, I.; Wang, Y.; Arora, S.; Atkinson, D.; Authur, R.; Chandu, K.; Cohan, A.; Dumas, J.; Elazar, Y.; Gu, Y.; Hessel, J.; Khot, T.; Merrill, W.; Morrison, J.; Muennighoff, N.; Naik, A.; Nam, C.; Peters, M. E.; Pyatkin, V.; Ravichander, A.; Schwenk, D.; Shah, S.; Smith, W.; Subramani, N.; Wortsman, M.; Dasigi, P.; Lambert, N.; Richardson, K.; Dodge, J.; Lo, K.; Soldaini, L.; Smith, N. A.; and Hajishirzi, H. 2024. OLMo: Accelerating the Science of Language Models. *Preprint*.

Grynbaum, M. M.; and Mac, R. 2023. The Times Sues OpenAI and Microsoft Over A.I. Use of Copyrighted Work.

Hendrycks, D.; Burns, C.; Basart, S.; Zou, A.; Mazeika, M.; Song, D.; and Steinhardt, J. 2021. Measuring Massive Multitask Language Understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*.

Jiang, A. Q.; Sablayrolles, A.; Mensch, A.; Bamford, C.; Chaplot, D. S.; de las Casas, D.; Bressand, F.; Lengyel, G.; Lample, G.; Saulnier, L.; Lavaud, L. R.; Lachaux, M.-A.; Stock, P.; Scao, T. L.; Lavril, T.; Wang, T.; Lacroix, T.; and Sayed, W. E. 2023. Mistral 7B. arXiv:2310.06825.

Lin, C.-Y. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, 74–81. Barcelona, Spain: Association for Computational Linguistics.

Liu, B.; Liu, Q.; and Stone, P. 2022. Continual learning and private unlearning. In *Conference on Lifelong Learning Agents*, 243–254. PMLR.

Maini, P.; Feng, Z.; Schwarzschild, A.; Lipton, Z. C.; and Kolter, J. Z. 2024. TOFU: A Task of Fictitious Unlearning for LLMs. arXiv:2401.06121.

Mattei, S. E.-D. 2023. Artists Are Suing Artificial Intelligence Companies and the Lawsuit Could Upend Legal Precedents Around Art. Pawelczyk, M.; Neel, S.; and Lakkaraju, H. 2024. In-Context Unlearning: Language Models as Few Shot Unlearners. In *ICML*.

Shi, W.; Lee, J.; Huang, Y.; Malladi, S.; Zhao, J.; Holtzman, A.; Liu, D.; Zettlemoyer, L.; Smith, N. A.; and Zhang, C. 2024. MUSE: Machine Unlearning Six-Way Evaluation for Language Models. arXiv:2407.06460.

Soldaini, L.; Kinney, R.; Bhagia, A.; Schwenk, D.; Atkinson, D.; Authur, R.; Bogin, B.; Chandu, K.; Dumas, J.; Elazar, Y.; Hofmann, V.; Jha, A. H.; Kumar, S.; Lucy, L.; Lyu, X.; Lambert, N.; Magnusson, I.; Morrison, J.; Muennighoff, N.; Naik, A.; Nam, C.; Peters, M. E.; Ravichander, A.; Richardson, K.; Shen, Z.; Strubell, E.; Subramani, N.; Tafjord, O.; Walsh, P.; Zettlemoyer, L.; Smith, N. A.; Hajishirzi, H.; Beltagy, I.; Groeneveld, D.; Dodge, J.; and Lo, K. 2024. Dolma: an Open Corpus of Three Trillion Tokens for Language Model Pretraining Research. *arXiv preprint*.

ssa. 2011. Social Security is Changing the Way SSNs are Issued. https://www.ssa.gov/kc/SSAFactSheet--IssuingSSNs. pdf. Accessed: 2024-10-07.

Triantafillou, E.; Pedregosa, F.; Hayes, J.; Kairouz, P.; Guyon, I.; Kurmanji, M.; Dziugaite, G. K.; Triantafillou, P.; Zhao, K.; Hosoya, L. S.; Junior, J. C. S. J.; Dumoulin, V.; Mitliagkas, I.; Escalera, S.; Wan, J.; Dane, S.; Demkin, M.; and Reade, W. 2023. NeurIPS 2023 - Machine Unlearning.

Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837.

Zhang, R.; Lin, L.; Bai, Y.; and Mei, S. 2024. Negative Preference Optimization: From Catastrophic Collapse to Effective Unlearning. arXiv:2404.05868.

Zhao, K.; Kurmanji, M.; Bărbulescu, G.-O.; Triantafillou, E.; and Triantafillou, P. 2024. What makes unlearning hard and what to do about it. arXiv:2406.01257.

Limitations and Future Work

Carlini et al. (2022) show that the risk of memorization increases with large model size. However, due to computational limitations and easy availability of large public LLMs, we only provide finetuned checkpoints for 1B and 7B OLMo, and defer release of larger models to future work. Moreover, licensing restrictions prevent us from releasing fine-tuned models based on few publicly available LLMs such as LLaMa (lla 2023).

We acknowledge that LLM-generated data can exhibit specific biases found in their training data set. We partially mitigate this by seeding the generation prompt with pre-sampled character and location names to ensure diversity in generated content. We also conducted manual evaluations of the generated creative content to ensure its quality.

Ethical Considerations

Task 2 deals with sensitive PII information which warrants careful considerations to avoid privacy leakage of individuals. We avoid this risk entirely by carefully designing the generation process so that it closely mimics real individuals, despite being generated synthetically. We also ensure all the tools used in generating our benchmark data are open sourced, thereby avoiding any licensing restrictions.

Long Form Example Generation

Model Input:

Create a short novel of at least 150 words. The novel should be from fantasy genre and set in the city of Atlantis. It should include following characters: Kyle, Stan, Kenny and Eric.

Model Output:

A short story about four friends lost in the magical realm below the oceans, known to ousiders as Atlantis. Kyle had always held a deep fascination for the deep blue ocean, and this naturally led him to take up a major in oceanic studies...

Short Form Example Generation

Model Input:

Create a biography for Jon Smith with date of birth: 1/2/1989, SSN: 900123456, phone number: 0987654321, email: jon_smith@me.com, home address: 10 Summertime Lane, New York City, NY, USA.

Model Output:

Jon Smith was both in New York City on the first of February in 1989,...

Question Generation Prompt

Model Input:

You are given a short story. First, find all the proper nouns in this story. If it does not contain a proper noun, say " I can't use this statement since it does not contain any proper nouns.". If it contains proper nouns, use your reasoning to create an unambiguous question, for which there would be *only* one answer. Give a concise answer (i.e . one word or phrase) which accurately answers the question. If you cannot create such an unambiguous question, say "I'm unable to create an unambiguous question for this story". Use the examples below for reference.
Examples:

Example #1
Example #2
Example #3
Example #4
Example #4

Here's the story: <input_story>. Generate a question with an unambiguous answer using this story.

Further details on Unlearning Algorithms

We review unlearning methods tested in this paper in the following.

• **Gradient Ascent**: This is a straightforward algorithm for model unlearning where we reverse the direction of model update by flipping the sign in gradient descent, in order to steer the model away from the sensitive model outputs in the forget set. While easy to implement, this approach has a significant drawback since the gradient ascent training objective is unbounded,

which can lead to model divergence with nonsensical outputs for all inputs. The loss term in this algorithm reverses sign of the standard training objective and is applied only on the forget set F as shown below.

 $-\mathcal{L}(F;\theta)$

• Gradient Difference (Liu, Liu, and Stone 2022): In this approach, we augment the gradient ascent objective applied on forget set, by adding a gradient descent objective on the retain set. By jointly optimizing on both sets, we steer the model away from regurgitating the sensitive information from the retain set, while ensuring it does not lose performance in the retain set. Despite being a promising alternative to Gradient Ascent, this quality of model performance on non-sensitive dataset depends on the size of the retain set used in model training, and can lead to poor generalization on new examples. The loss term jointly increases the likelihood of generating responses in the retain set *R* while reducing the likelihood of generating *F*, as shown below.

$$-\mathcal{L}(F;\theta) + \mathcal{L}(R;\theta)$$

- **KL Divergence** (Maini et al. 2024) Similar to Gradient Difference, in this baseline, we augment the gradient ascent objective with a Kullback-Leibler Divergence term to ensure the model does not deviate too far from the original model.
- Negative Preference Optimization (Zhang et al. 2024): This baseline uses a modified version of the Direct Preference Optimization objective, adapted to remove the sensitive information from the forget set.

More details on related work

Given the nascent stage of unlearning research in LLMs, few prior works exist which address the task of robustly evaluating the success of unlearning. Triantafillou et al. (2023) presented a new challenge task in which the goal was to to unlearn information contained in select images within the task of image based age prediction. While successful, the specific task addressed in this challenge was narrow, focusing only on image based age prediction - a classification problem with 10 classes with limited applicability in the unbounded text generation task of large language models. But the growing interest in LLMs and their tendency to regurgitate sensitive or private information necessitates a distinct and focused evaluation benchmark.

Maini et al. (2024) released a new evaluation framework named TOFU which partially addressed this task of evaluating LLM unlearning algorithms. Their framework was evaluated on question answering task applied on biographies of synthetically created fake authors. They train target models on this synthetic data and evaluate the ability of unlearning algorithms to forget a portion of this synthetic dataset. While being a promising first step, this work has a few key limitations: unlearning the targeted information required for the QA task does is unlikely to cause loss of any other substantial information, specially linguistic attributes such as grammar. Further, this work leverages GPT4 to generate the synthetic content, which may have downstream licensing implications owing to GPT4's proprietary license.

More recently, Shi et al. (2024) released a benchmark named MUSE which evaluated model unlearning using real data set for containing news documents and Harry Potter book chapters. This benchmark released detailed evaluation metrics to robustly evaluate the unlearning algorithms. However since it only leverages real data set the benchmark does not provide a clean test bed to evaluate model performance. Specifically, the information contained in the unlearn documents may also appear in other disjoint training documents, limiting the effectiveness of unlearning. While the TOFU benchmark mentioned before avoids this by only using synthetic documents, the data set coverage is rather limited (it only containts biographic information). The benchmark developed in this work addresses both these shortcomings together and presents a single holistic testbed to evaluate model unlearning in LLMs.