

The Landscape of Memorization in LLMs: Mechanisms, Measurement, and Mitigation

Alexander Xiong¹, Xuandong Zhao¹, Aneesh Pappu², Dawn Song¹

¹University of California, Berkeley

²Google DeepMind

{alexixiong,xuandongzhao,dawnsong}@berkeley.edu, aneeshpappu@google.com

Abstract

Large Language Models (LLMs) have demonstrated remarkable capabilities across a wide range of tasks, yet they also exhibit memorization of their training data. This phenomenon raises critical questions about model behavior, privacy risks, and the boundary between learning and memorization. Addressing these concerns, this paper synthesizes recent studies and investigates the landscape of memorization, the factors influencing it, and methods for its detection and mitigation. We explore key drivers, including training data duplication, training dynamics, and fine-tuning procedures that influence data memorization. In addition, we examine methodologies such as prefix-based extraction, membership inference, and adversarial prompting, assessing their effectiveness in detecting and measuring memorized content. Beyond technical analysis, we also explore the broader implications of memorization, including the legal and ethical implications. Finally, we discuss mitigation strategies, including data cleaning, differential privacy, and post-training unlearning, while highlighting open challenges in balancing the minimization of harmful memorization with utility. This paper provides a comprehensive overview of the current state of research on LLM memorization across technical, privacy, and performance dimensions, identifying critical directions for future work.

1 Introduction

Throughout the past few years, we have observed great strides in the capabilities of LLMs driven by changes in model architecture, training methodologies, and computational resources [Radford et al., 2018, Brown et al., 2020, Chowdhery et al., 2023, Naveed et al., 2023, Touvron et al., 2023, Wei et al., 2023]. These advances have significantly improved natural language understanding, generation, and reasoning abilities, enabling these models to perform increasingly complex tasks across diverse domains [Alowais et al., 2023, Khurana et al., 2023, Minaee et al., 2025]. However, in addition to their strong capabilities, LLMs continue to manifest critical limitations when evaluated on criteria such as privacy and security [Chang et al., 2024, Yao et al., 2024a, Das et al., 2025].

With the advent of data scale used to train LLMs, the risks of privacy leakages through memorization have dramatically increased, as models trained on massive datasets containing potentially sensitive information may inadvertently reproduce verbatim passages from training data when prompted with specific triggers, compromising individual privacy and confidentiality without explicit consent [Dwivedi et al., 2023, Smith et al., 2023, Abdali et al., 2024]. Thus, data memorization emerges as a critical vulnerability in LLMs, presenting significant privacy risks, including potential exposure of personally identifiable information (PII), copyright violations, and unauthorized reproduction of

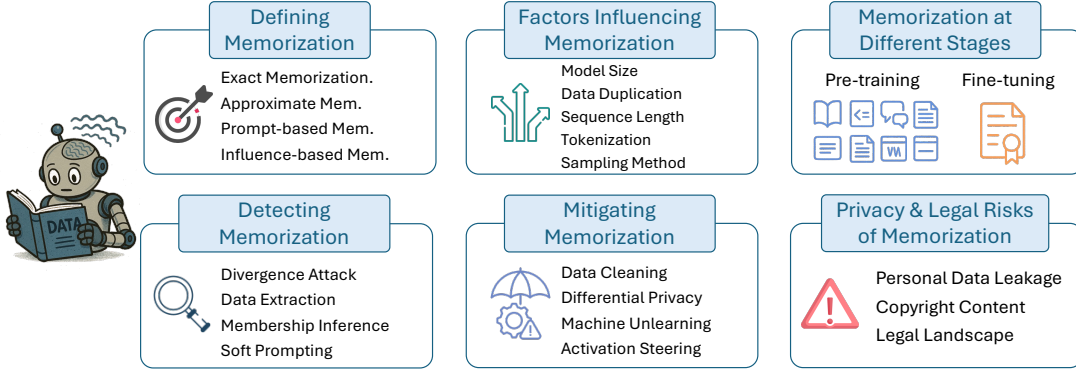


Figure 1: Taxonomy of understanding memorization in LLMs.

sensitive content [Carlini et al., 2021, Cremer et al., 2022, Clusmann et al., 2023, Karamolegkou et al., 2023, Song et al., 2023, Freeman et al., 2024, Hua et al., 2024, Morris et al., 2025].

This investigation explores memorization mechanisms in LLMs, examining contributing factors, detection methodologies, measurement approaches, and mitigation techniques. Unlike previous surveys on security and privacy [Siau and Wang, 2020, Fui-Hoon Nah et al., 2023, Guo et al., 2023, Hadi et al., 2023, Karabacak and Margetis, 2023, Min et al., 2023, Zhao et al., 2023, Zhu et al., 2023, Chang et al., 2024, Liu et al., 2024a, Myers et al., 2024, Raiaan et al., 2024, Stahl and Eke, 2024, Zhao et al., 2024], we present a focused examination of memorization phenomena in LLMs, dissecting the fundamental components contributing to and mitigating unintended data retention [Bender et al., 2021, Zhao et al., 2022a, Hartmann et al., 2023, Satvaty et al., 2024]. This paper gives an in-depth review of the existing research on LLM memorization and its outcomes. Opposed to the survey by Satvaty et al. [2024], our analysis examines the technical underpinnings, evaluation methods, and privacy implications of memorization while providing concrete research directions. Figure 1 provides a visual taxonomy that maps and categorizes memorization, detection & mitigation of memorization, and corresponding risks.

2 Defining memorization

Agnostic to LLMs, memorization occurs when a model can reproduce specific training sequences when provided with appropriate contextual inputs [Carlini et al., 2023, Cooper et al., 2023, Schwarzschild et al., 2024]. Currently, several definitions are used to quantify memorization in LLMs, including memorization defined by string matching (exact or approximate) or by loss differences.

2.1 Exact memorization

Verbatim memorization refers to a model’s exact reproduction of training data, often arising due to the high duplication of specific examples or overfitting [Carlini et al., 2021]. This phenomenon occurs when an LLM, instead of generating novel responses or synthesizing information, simply retrieves and reproduces specific text snippets word-for-word from the data it was exposed to.

Perfect memorization describes a model that has exactly recorded its training data, assigning a probability only to inputs it has previously seen [Kandpal et al., 2022]. When such a model generates an output, the process is mathematically equivalent to randomly selecting examples from its training dataset with a uniform probability.

Eidetic memorization [Carlini et al., 2021] occurs when a prompt p causes the model to reproduce a verbatim string s from training data. A stricter variant, **k -eidetic memorization** [Carlini et al., 2021], occurs when s appears in no more than k training examples, yet can still be extracted given the right prompt.

Discoverable memorization. A string s is discoverably memorized [Carlini et al., 2023, Nasr et al., 2023] when, given a training example composed of prefix p and suffix s , a LLM exactly reproduces the continuation s that followed p in the training data.

2.2 Approximate memorization

Approximate/Paraphrased memorization [Ippolito et al., 2023] is when LLMs generate outputs similar to training data in content, structure, or phrasing without exact replication. This differs from exact copying through variations, paraphrasing, or partial overlaps. The authors detect memorization by calculating edit distance between a generation and a target string, normalized by length, and choosing a threshold that determines when a sequence is deemed an approximate match.

2.3 Prompt-based memorization

Extractable memorization occurs when, without access to the training data, there exists a constructable prompt that invokes the model to generate an exact example from its training set [Nasr et al., 2023].

k -extractable memorization represents a stricter form of extractable memorization. A suffix is k -extractable [Biderman et al., 2023]. when, given only its corresponding k -token prefix, the model reproduces the entire suffix verbatim. Unlike verbatim memorization, which manifests as direct replication, k -extractable memorization captures a model’s ability to retrieve and complete specific training sequences when prompted with only partial context.

(n, p) -discoverable extraction [Hayes et al., 2025] formalizes the likelihood of retrieving a memorized string via repeated sampling. A string is (n, p) -discoverable if it appears in at least one of n completions with probability $\geq p$.

2.4 Influence-based memorization

Counterfactual memorization [Feldman and Zhang, 2020, Zhang et al., 2023] quantifies how much a model’s predictions change based on whether a model is trained on a particular training example. For example, this can be instantiated as how much the loss changes on a datapoint when seen by a model during training versus not. Pappu et al. [2024] instantiate this definition by classifying datapoints as memorized when both a model trained on the datapoint achieves an approximate match under prefix-based decoding and a model not trained on that datapoint does not achieve an approximate match. As it is computationally prohibitive to retrain a model per datapoint-exclusion, Zhang et al. [2023] and Pappu et al. [2024] train a smaller number of models that exclude entire subsets of data used to test for memorization.

3 Factors influencing memorization

With the increasing scale and capability of LLMs, a range of training and inference-time factors have been identified as influencing memorization behavior. This section systematically analyzes the

primary variables that govern memorization.

Model parameter size. Model size strongly correlates with increased memorization in LLMs, with larger models demonstrating both greater capacity to retain training data and increased vulnerability to extraction attacks. [Carlini et al. \[2021\]](#) first introduced the relationship between model size and memorization, observing that memorization scales log-linearly with model size. Subsequent research has consistently reinforced this finding across multiple dimensions of memorization. For instance, [Tirumala et al. \[2022\]](#) observed that larger LLMs not only memorize more content but do so more rapidly during training, a phenomenon that cannot be fully explained by conventional overfitting or hyper-parameter tuning. This enhanced memorization capability directly impacts privacy concerns, as [Li et al. \[2024\]](#) demonstrated that extraction attacks become significantly more effective against larger LLMs due to their advanced memorization capabilities. The scaling effect has been further corroborated by [\[Lee et al., 2022, Tirumala et al., 2022, Carlini et al., 2023, Freeman et al., 2024, Kiyomaru et al., 2024, Li et al., 2024, Hayes et al., 2025, Morris et al., 2025\]](#).

Training data duplication. Duplicate examples in training data skew LLMs toward overrepresented content, diminishing output diversity and increasing verbatim reproduction of training material [\[Carlini et al., 2021\]](#). As demonstrated in [Lee et al. \[2022\]](#), de-duplication substantially reduces memorization, with models trained on original data exhibiting a tenfold increase in memorized token generation compared to those trained on datasets, where exact substrings and near-duplicates were removed. Exploring further, [Kandpal et al. \[2022\]](#) identified a superlinear relationship between training data duplication and memorization in LLMs, where rarely duplicated training samples are seldom memorized. However, models generate duplicated sequences at rates significantly lower than perfect memorization would predict.

Sequence length. Longer sequences increase memorization in LLMs. [Carlini et al. \[2023\]](#) found that memorization increases logarithmically with sequence length, with the verbatim reproduction probability rising by orders of magnitude as sequences extend from 50 to 950 tokens. Similarly, extraction methodologies, including continuous soft prompting and dynamic soft prompting techniques, demonstrate a consistent pattern wherein the volume of extracted data grows proportionally with increasing prefix token sizes [\[Wang et al., 2024\]](#).

Tokenization. Models trained with larger Byte Pair Encoding (BPE) vocabularies were found to memorize significantly more sequences from their training data. [Kharitonov et al. \[2021\]](#) studied how different vocabulary sizes, produced via BPE, affect a transformer LLMs’ tendency to memorize sequences. Memorization was particularly strong for named entities, URLs, and uncommon phrases, often becoming single tokens under larger BPE settings.

Sampling Methods. Sampling methods significantly impact the extraction of memorized content from LLMs, with stochastic approaches consistently outperforming greedy decoding in revealing memorized training data. While [Carlini et al. \[2021\]](#) initially established that greedy decoding can reveal memorized sequences but misses others due to its low diversity, subsequent research has demonstrated the superior effectiveness of other sampling methods. [Yu et al. \[2023\]](#) showed that optimizing sampling parameters like top-k, nucleus sampling, and temperature can substantially increase memorized data extraction, in some cases doubling previous baselines. This finding was reinforced by [Tiwari and Suh \[2025\]](#), who discovered that randomized decoding nearly doubles leakage risk compared to greedy decoding, contradicting earlier assumptions. [Hayes et al. \[2025\]](#) provided a theoretical foundation for these observations by introducing a probabilistic framework that demonstrates how repeated sampling with varied decoding parameters can expose memorization hidden under greedy approaches. Given that no single decoding method minimizes leakage across

all scenarios, [Tiwari and Suh \[2025\]](#) emphasized the importance of assessing memorization under diverse sampling strategies.

While data duplication has been broadly associated with increased memorization, it remains uncertain whether this effect systematically varies across different data modalities (e.g., code, prose, etc.). Additionally, although inference hyper-parameters, such as decoding methods, are known to impact memorization, existing techniques do not reliably prevent the extraction of memorized content. This underscores the need for decoding strategies explicitly designed to minimize memorization. In contrast, the influence of training hyper-parameters on memorization remains poorly understood, particularly with regards to the role of overfitting through learning rate and regularization mechanisms. Moreover, the temporal dynamics of memorization in the presence of data duplication remain underexplored. Specifically, it is unclear how the timing of when duplicated examples are encountered during training affects memorization likelihood.

Open Questions

1. How do data duplication and data type affect memorization extent?
2. What is required to design a decoding method that minimizes leakage?
3. What factors cause larger LLMs to memorize more?
4. How can we distinguish useful generalization from memorization in LLMs?
5. How do training hyperparameters influence memorization?
6. Does the temporal order of data presentation during training affect memorization likelihood?

4 Memorization at different stages

Although memorization has often been attributed to overfitting, this section explores how it is more deeply shaped by pre-training dynamics and fine-tuning strategies in LLMs.

Pre-training dynamics introduce systematic biases that influence which training examples are retained. Under standard non-deterministic training regimes, characterized by data shuffling, dropout, and stochastic optimization, models tend to forget examples seen early in training unless they are revisited, as shown in [Jagielski et al. \[2023\]](#), [Kiyomaru et al. \[2024\]](#) and [Leybzon and Kervadec \[2024\]](#). They hypothesize that this forgetting is not simply due to limited exposure but rather the result of parameter drift, whereby updates driven by later data overwrite earlier representations. Consequently, memorization is biased toward examples encountered in the latter stages of training.

This effect becomes even more pronounced as models approach the end of training. [Huang et al. \[2024\]](#) showed that later-stage checkpoints are more susceptible to memorizing even rare or out-of-distribution content, likely reflecting increased model capacity and representational flexibility at those stages. Extending this perspective, [Biderman et al. \[2023\]](#) demonstrated that memorization follows predictable scaling laws. As model size and training duration increase, specific sequences predictably transition from unmemorized to memorized. This behavior reinforces that memorization is not an artifact of overfitting but an emergent property of scale and optimization dynamics.

Supervised fine-tuning influences memorization behavior, shaping both the extent and nature of information retention and exposure. [Mireshghallah et al. \[2022\]](#) provided a comparative analysis of fine-tuning approaches, demonstrating that head-only fine-tuning presents the highest risk of memorization, likely due to overfitting. In contrast, adapter-based fine-tuning, when constraining

parameter updates [Zeng et al., 2024], reduces memorization using parameter-efficient methods, demonstrating fine-tuning can be effective for limiting privacy risks. They link differences in memorization from models fine-tuned for tasks such as summarization vs. question answering to the attention dynamics of each task, with narrower attention patterns correlating with higher memorization. From an adversarial perspective, Chen et al. [2024] demonstrated that fine-tuning can also be leveraged to extract memorized pretraining data. Using the Janus Interface, they show that targeted fine-tuning amplifies a model’s capacity to leak sensitive information, thus framing fine-tuning as a mechanism that can reactivate latent memorization vulnerabilities. Together, these studies underscore that fine-tuning critically shapes memorization patterns.

Reinforcement learning as post-training. Post-training paradigms such as Reinforcement Learning from Human Feedback (RLHF) [Ouyang et al., 2022], Reinforcement Learning with Verifiable Rewards (RLVR) [Lambert et al., 2024], and Reinforcement Learning from Internal Feedback (RLIF) [Zhao et al., 2025] leverage reinforcement learning to train LLMs from various feedback sources. Limited research exists on memorization propagation through reinforcement learning stages. In the context of code generation, Pappu et al. [2024] found that data memorized during fine-tuning persists with high frequency in post-RLHF models, while finding minimal evidence that reward model data or reinforcement learning data becomes memorized. Critical questions remain regarding the predictability of memorization persistence between fine-tuning and reinforcement learning based on data attributes.

Distillation is another common technique used in modern ML pipelines [Hinton et al., 2015, Zhao et al., 2022b], where a small ‘student’ model is trained on data produced by a larger ‘teacher’ model (typically logits). Chaudhari et al. [2025] show that bias injected adversarially into teacher models can be amplified in student models via distillation. This naturally suggests that memorization can propagate between teacher and student models, but this has not been formally analyzed.

Characterizing memorization across training phases requires understanding the mechanisms by which information is encoded, retained, or discarded. Although memorization tends to intensify in later training stages, the degree to which it can be attenuated or reversed remains an open question. With the increasing prevalence of fine-tuned models, it is also unclear how established memorization scaling laws extend to domain-adapted settings. Moreover, the relative contributions of training dynamics versus fine-tuning objectives to overall memorization behavior are not yet well understood.

Open Questions

1. How can we reduce parameter drift to retain early training knowledge?
2. Do memorization scaling laws apply to domain-specific fine-tuned models?
3. How can we quantify and attribute memorization to pre-training vs. fine-tuning?
4. Are there identifiable factors that predict what type of data will persist as memorized between pre-training and post-training?
5. How does memorization persist (or not) from teacher to student models during distillation?

5 Detecting memorization

A wide range of techniques has been developed to detect memorization in LLMs. This section categorizes and examines these methods in depth to offer a comprehensive understanding of how memorization is identified.

Divergence attack. Nasr et al. [2023] proposed a divergence attack, a prompt-based extraction method that circumvents alignment defenses in instruction-tuned LLMs, coaxing the model toward reverting to pre-aligned behavior. This divergence significantly increases the likelihood of emitting memorized training data, achieving up to $150\times$ more verbatim sequences compared to typical user queries. Nasr et al. [2023] hypothesized that this vulnerability arises because prompts induce decoding analogous to an end-of-text token during pretraining, a context in which LLMs are known to favor high-likelihood and often memorized continuations. The attack exposes persistent memorization that may have remained latent.

Prefix-based data extraction attack works by querying a model with the initial segment of a memorized sequence and observing whether it completes the rest (verbatim or approximate, depending on the specific definition used). First proposed by Carlini et al. [2021], they demonstrated this technique on GPT-2, successfully recovering hundreds of memorized examples by prompting the model with partial strings from the training set. Carlini et al. [2023] showed that longer prefixes increased the likelihood of verbatim completions by reducing ambiguity and aligning the model more strongly with memorized content. Building on Carlini et al. [2023], Li et al. [2024] introduced the LLM-PBE toolkit, a benchmark for evaluating privacy risks via prefix-based extraction. Li et al. [2024]’s findings reinforced that structured prefixes, such as email headers or document beginnings, are potent at eliciting memorized sequences.

An adversarial variant of this attack was explored by Kassem et al. [2024], who used LLMs to generate candidate prefixes likely to elicit private data from a target model. Kassem et al. [2024] observed that instruction-tuned models can leak pretraining data even when the prompts diverge from the original training distribution, suggesting that memorization in LLMs may be more pervasive than what is revealed by prefix attacks. An extension of prefix attacks is soft prompting, which is explored in further detail below.

Membership inference attack (MIA) [Shokri et al., 2017] has emerged as a key diagnostic tool for detecting memorization in LLMs, offering insight into whether specific data points may have been memorized during training. To circumvent the computational overhead of training shadow models, recent work has focused on black-box strategies that exploit the model’s output probabilities or perplexity for input sequences [Duan et al., 2024].

A number of MIA techniques rely on loss- or likelihood-based metrics. Yeom et al. [2018] used the raw loss on a target input, assuming seen examples have lower loss. To mitigate confounding effects from input difficulty, reference-based calibration [Carlini et al., 2021] subtracts the loss from a separate reference model, aiming to isolate differences due to training exposure. Another method used in Carlini et al. [2021] is zlib entropy, which normalizes loss by the zlib-compressed size of the input sequence, using compression as a proxy for complexity. Moving beyond single-point estimates, the neighborhood attack [Mattern et al., 2023] examines local changes in the loss landscape. By generating nearby perturbations of the target input and comparing their average loss, it identifies instances where the target appears memorized. Similarly, the min- $k\%$ prob method [Shi et al., 2023] focuses on the subset of tokens with the highest loss—under the assumption that even the most uncertain tokens in a memorized sequence will still be predicted with confidence.

Yet, despite their utility in highlighting patterns of memorization, MIAs fall short when used as per-instance indicators of training data inclusion. As Zhang et al. [2025] and Duan et al. [2024] argued, MIAs lack a well-calibrated null model since one cannot feasibly train an identical model without the target input. This makes it impossible to meaningfully estimate false positive rates, undermining the statistical soundness of individual predictions. This limitation is underscored by

the critique in Maini and Suri [2024], which challenges the claims made in Zhang et al. [2024]. The latter introduces the PatentMIA benchmark and a divergence-based calibration method for MIAs, reporting improved reliability in detecting whether a sequence appeared in pretraining. However, Maini and Suri [2024] showed these gains are likely artifacts of distributional shifts introduced by the benchmark itself, which MIAs can exploit without truly resolving the underlying calibration problem. In light of these concerns, Zhang et al. [2025] recommended that MIAs be re-framed as tools for aggregate-level privacy auditing rather than as evidence of individual training data exposure. For stronger guarantees, they advocate for alternatives, such as data extraction attacks or canary-based MIAs, which offer more direct insights into memorization and leakage.

Soft prompting techniques have emerged as powerful tools, enabling both the extraction and suppression of memorized content. By leveraging learned embeddings pre-pended to model inputs, continuous soft prompts can explicitly influence what a model reveals from its training data. For instance, Ozdayi et al. [2023] demonstrated that fixed-length continuous prompts could be trained to either amplify or diminish memorization, as measured by the extraction rates of secret sequences. Their results showed that attack prompts increased memorization leakage by up to 9.3%, while suppress prompts decreased extraction by up to 97.7%. Building on Ozdayi et al. [2023], Kim et al. [2023] introduced ProPILE, a privacy auditing framework that uses soft prompt tuning in a white-box setting to extract PII memorized by LLMs. Their learned prompts were transferable across models, suggesting consistent memorization patterns that can be systematically exploited.

Dynamic soft prompting techniques extend the static nature of continuous prompts by conditioning on the input context. Wang et al. [2024] proposed prefix-conditioned prompting, where a prompt generator produces tailored soft prompts based on the input prefix. This dynamic approach allows LLMs to adapt to subtle input variations and surface context-dependent memorized completions. Wang et al. [2024] revealed that dynamic prompts significantly improved the discoverable memorization rate, surpassing both static prompt and no-prompt baselines. These results underscore the role of context-sensitive prompting in uncovering latent memorized content and highlight the broader utility of soft prompting in analyzing memorization in LLMs.

Reasoning. Much of recent literature has been focused on training Large Reasoning Models (LRMs), models trained to reason before producing a final answer, as opposed to immediately producing an answer in their first token. To date, most memorization literature has been focused on quantifying memorization of specific datapoints by way of measuring whether the model can reproduce the training example. As the field moves towards reasoning models, a definition of what it means for a model to memorize a reasoning pattern may be useful for guiding reasoning research and ensuring reasoning models generalize. Huang et al. [2025] takes a step towards this by measuring whether models trained to reason on math are able to solve slightly perturbed versions of problems in their training set. The authors find that models fail to solve problems that constitute ‘hard’ perturbations, where the problem looks superficially similar to a problem in the training set but requires a different reasoning strategy. Creating benchmarks to measure memorization of such superficial reasoning patterns remains an open and useful direction for LRM research.

From the memorization detection strategies explored above, except for the divergence attack, adversaries either imbue or infer/know the training data intended for extraction. Therefore, it is critical to develop methods to reliably detect memorized sequences in language models without requiring access to the original training data. Furthermore, these detection methods all examine direct memorization from training data. With more powerful LLMs, the boundary between memorization and generalization is fuzzier, propelling the need to understand if memorization can be detected from patterns using reasoning.

Open Questions

1. How can we reliably detect memorization without training data access?
2. Can LLMs infer training data from context rather than direct memorization?
3. How can we detect memorization of implicit patterns and sensitive relationships?
4. How can we build benchmarks to measure whether reasoning models memorize superficial reasoning patterns or learn general reasoning principles?

6 Mitigating memorization

Mitigation techniques for memorization span various approaches, from data handling to model training and auditing. This section outlines the main categories, explores key strategies within each, and highlights their differences in effectiveness and trade-offs.

6.1 Training-time interventions

Data cleaning plays a crucial role in mitigating memorization by limiting an LLM’s exposure to overrepresented sequences. As demonstrated by [Carlini et al. \[2021\]](#), [Lee et al. \[2022\]](#), and [Kandpal et al. \[2022\]](#), de-duplication minimizes the over-representation of repeated sequences. Removing these duplicates serves as a form of implicit regularization, reducing the likelihood that the LLM overfits.

From an optimization perspective, LLMs can minimize loss on rare training examples by memorizing. This makes PII-scrubbing effective: by removing these sequences during training, memorization is minimized [[Carlini et al., 2021](#), [Jagielski et al., 2023](#)].

Differential privacy (DP) provides robust protection against membership inference and extraction attacks, even in the presence of adversaries with auxiliary knowledge [[Dwork, 2006](#), [Dwork et al., 2014](#), [Shokri et al., 2017](#), [Carlini et al., 2021](#)]. Differentially Private Stochastic Gradient Descent (DP-SGD) [[Abadi et al., 2016](#)] guarantees privacy by computing per-example gradients, clipping them to a fixed norm to bound individual influence, and injecting calibrated Gaussian noise into the aggregated updates. A privacy accountant tracks cumulative privacy loss to enforce the global (ϵ, δ) budget.

[Li et al. \[2021\]](#) demonstrated that LLMs can act as strong DP learners. Their evaluation reveals that pretrained LLMs are highly resilient to the noise introduced by DP-SGD, particularly when fine-tuned on top of general-purpose representations. DP-trained models approach the performance of non-private baselines, while offering provable protection against memorization. Empirical validation through canary insertion and MIAs confirm the effectiveness of DP in preventing leakage of sensitive sequences. To further reduce memorization, [Zhao et al. \[2022a\]](#) introduced Confidentially Redacted Training (CRT), a method that combines de-duplication and redaction with DP-SGD to train LLMs while avoiding the retention of sensitive content. CRT builds on DP concepts to introduce randomized training interventions that prevent unintended memorization. The authors show that CRT, when combined with DP-SGD, provides strong privacy protections without compromising model performance, as indicated by competitive perplexity scores.

However, the use of DP in LLMs often presents trade-offs between privacy, utility, and computational efficiency. Strict privacy budgets often lead to performance drops, and significant computational overhead makes it challenging to scale to LLMs. To address these challenges, recent work has

examined parameter-efficient fine-tuning (PEFT) strategies for training under DP. PEFT approaches, such as LoRA [Hu et al., 2022], allow finetuning of models with addition of few parameters relative to base models via low-rank decompositions of additional trainable parameters. It is hypothesized that finetuning PEFT models with DP mitigates losses of utility due to requiring less noise for privatization due to training with fewer parameters [Liu et al., 2024b, Yu et al., 2022, 2021]. Ma et al. [2024] explored the integration of adapter-based fine-tuning methods into the DP framework. Their findings reveal that DP PEFT reduces memorization and achieves comparable or superior downstream task performance relative to full-model DP-SGD, under strict privacy budgets. Ma et al. [2024] hypothesized that the limited parameter updates in PEFT may concentrate DP noise in a narrow subset of the model, potentially weakening privacy protection. While most DP techniques for LLMs target record-level privacy, many applications require user-level guarantees. Chua et al. [2024] addressed this gap by proposing several techniques to enforce DP at user-level. To evaluate effectiveness, Chua et al. [2024] used canary insertion attacks, embedding unique sequences into user data to test for memorization. Results show that user-level DP reduces memorization, with much lower canary extraction rates than record-level DP.

DP in LLMs continue to present trade-offs between privacy, utility, and computational efficiency. Strict privacy budgets often lead to performance drops, and significant computational overhead makes it challenging to scale to LLMs. Thus, the adoption of DP for mitigating memorization in LLMs remains limited.

6.2 Post-training-time interventions

Machine unlearning aims to remove the influence of certain training examples so that the model’s behavior is as if those examples were never seen [Cooper et al., 2024, Liu, 2024]. Yao et al. [2024b] evaluates various model-editing unlearning strategies (e.g., gradient ascent [Golatkhar et al., 2020, Jia et al., 2023, Jang et al., 2022], negative re-labeling [Golatkhar et al., 2020], adversarial sampling [Cha et al., 2024]) demonstrating that approximate unlearning methods can be over $10^5 \times$ more computationally efficient than retraining a model from scratch. However, unlike DP, there is no formal guarantee, thereby leaving a risk that memorization persists.

ParaPO, introduced by Chen et al. [2025], is a post-training strategy to decrease memorization of data memorized in the pretraining corpus. ParaPO decreases unintended memorization by first assessing which data from the pretraining corpus has been memorized (by searching for sequences that, when a prefix is used to prompt an LM, an LM decodes a near exact match to the suffix of the sequence), and then creating preference pairs by using a separate LLM to summarize the memorized datapoint. The model is then post-trained via DPO [Rafailov et al., 2024] on pairs of (memorized sequence, summarization), with the summarization marked as the preferred response. Chen et al. [2025] show that this approach decreases unintended memorization while preserving verbatim recall of desired sequences (e.g. direct quotations) However, the authors note that this approach slightly decreases utility on math, knowledge and reasoning benchmarks.

Machine unlearning and ParaPO constitute post-training interventions for reducing memorization. ParaPO results in slight degradation of utility, raising the question of what extent memorization is required for utility and generalization. Additionally, a future research direction may be to what extent this tradeoff can be tuned via incorporation of “memorization” reward models that are optimized in conjunction with utility-oriented reward models, in the context of RLHF or RLVR.

6.3 Inference-time interventions

MemFree decoding is a strategy introduced by Ippolito et al. [2023] to filter out memorized sequences during generation. This approach contrasts with decoding by introducing a privacy filter into the generation loop applied post hoc at generation and acts as a wrapper around any pre-trained model. Ippolito et al. [2023] leveraged a bloom filter to represent the set of all n-grams in the training set to identify memorized content in real time. However, near-identical n-grams are undetectable as minor modifications to a token sequence evade detection. Another limitation is that access to the LLM’s n-gram training data is required.

Activation steering addresses memorization in LLMs by manipulating internal activations during inference [Turner et al., 2023]. By injecting targeted perturbations into an LLM’s hidden states, activation steering can suppress or redirect the generation of memorized sequences while preserving overall model reasoning capabilities [Suri et al., 2025]. Suri et al. [2025] demonstrated that using a sparse autoencoder to identify activation patterns linked to memorized passages enables the construction of steering directions that reduced memorization by up to 60%, with minimal performance degradation. However, they note that steering effectiveness depends heavily on precise layer selection and steering strength.

Understanding where memorization arises in LLMs is crucial to designing effective steering interventions. Stoehr et al. [2024] conducted a detailed mechanistic analysis of paragraph-length memorized sequences and found that memorization is often localized to specific LLM components. They showed that fine-tuning only a subset of the high-gradient weights was sufficient to erase memorized passages. Additionally, they identified a single attention head in an early layer that reliably activates in response to rare token combinations seen during training, triggering verbatim recall.

Building on localization, Chang et al. [2023] evaluated multiple methods for localizing memorized content by introducing two diagnostic benchmarks: injection, where known memorized sequences are inserted via fine-tuning a small set of weights, and deletion, which tests whether removing specific neurons erases a naturally memorized output. They compared several techniques, including activation-based heuristics [Geva et al., 2022], integrated gradients [Dai et al., 2021], and pruning-based methods like Slimming and Hard Concrete [Chang et al., 2023]. Their results showed that pruning-based approaches were most effective: for example, Hard Concrete was able to identify fewer than 0.5% of neurons whose removal led to a 60% drop in memorization accuracy. Both Chang et al. [2023] and Stoehr et al. [2024] highlighted a key challenge: the neurons involved in one memory often contribute to others, complicating memory-specific interventions due to the risk of collateral forgetting.

By leveraging insights from model localization and sparse representation learning, such techniques can target memorized content with increasing precision, offering a flexible and minimally invasive tool for minimizing memorization.

While data cleaning should be a fundamental part of every LLM training pipeline to reduce memorization, other mitigation strategies either face practicality and scalability trade-offs or lack formal guarantees of effectiveness. Memorization mitigation is typically approached from two perspectives: (1) addressing all training data, and (2) targeting specific subsets of information intended to be forgotten. Given that LLMs often indiscriminately memorize data, mitigation efforts should prioritize reducing the retention of sensitive or harmful content without compromising knowledge essential for factual recall and reasoning. From the first perspective, ensuring robust

privacy protection requires evaluating the viability of scaling up differentially private (DP) training and understanding the optimal trade-off between utility and cost. Additionally, as PEFT methods grow in popularity, there is a clear need to refine DP strategies to suit these approaches and mitigate some of the overhead associated with full-model DP training. For targeted mitigation, activation steering shows promise in localizing specific knowledge representations. However, further work is needed to formalize layer selection and parameter tuning to improve our understanding and control of memorization in LLMs.

6.4 Connections to other undesirable behaviors

It remains open whether mitigating memorization also mitigates other undesirable behaviors of LLMs, such as hallucination, racial bias, and toxic content. For instance, does reducing memorization reduce hallucination on concepts related to memorized data? Is it possible that reducing the memorization of datapoints that include harmful racial stereotypes reduces the rates at which large models produce biased or toxic content? The authors to date have not seen cross-cutting analysis examining whether reducing memorization can reduce other harmful behaviors in-tandem.

Open Questions

1. How can LLMs be designed or trained for selective memorization?
2. Shall we scale up DP training, and at what cost and utility trade-off?
3. How can DP improve parameter-efficient fine-tuning to reduce memorization?
4. How can we optimize activation steering for memorization removal?
5. How can post-training methods reduce memorization while preserving utility?
6. Can memorization mitigation methods be incorporated online during post-training, for instance, by using memorization-detecting reward models during RLHF?
7. When does memorization contribute to versus harm utility/generalization?
8. Does mitigating memorization affect other undesirable behaviors, such as hallucination or generating toxic stereotypes?

7 Privacy & legal risks of memorization

Memorization in LLMs poses serious security and privacy risks: memorized text can be directly leaked, undermining data confidentiality [Smith et al., 2023]. We discuss three major impact areas caused by memorization: personal data leakage, exposure of copyrighted or proprietary content, and broader legal and public consequences.

Personal data leakage. As explored in Li et al. [2024], Carlini et al. [2021], Jagielski et al. [2023], and Zhou et al. [2024], a significant risk for LLMs is leaking sensitive personal information. In Panda et al. [2024], researchers demonstrated a targeted neural phishing attack that achieved up to a 50% success rate in tricking a model into revealing PII during fine-tuning when injecting poisons during pretraining. Even when training data is intended to remain private, malicious attackers with the correct model access level can siphon memorized secrets [Panda et al., 2024]. If an LLM reveals any PII, this could result in identity theft or other harms. Furthermore, for industries such as healthcare or customer service, where customers expect confidentiality, this would place model operators in the crosshairs of current privacy regulations.

Copyright or proprietary content. Beyond PII, memorization in LLMs may also cause the reproduction of copyrighted or proprietary material such as books, code, etc. This turns the

LLM into a potential conduit for the unauthorized distribution of protected content. As observed by [Carlini et al. \[2023\]](#), models may produce verbatim memorization from their training data. While we can only explicitly confirm what training data is used for open-source LLMs [[Gao et al., 2020](#), [Soldaini et al., 2024](#)], there do exist copyrighted books in EleutherAI’s PILE dataset, specifically, copyrighted books in its Books3 section. This is observed in [Henderson et al. \[2023\]](#) and [Zhang et al. \[2022\]](#), where open models can regurgitate the first few pages of the Harry Potter books and “Oh the places you’ll go!” by Dr. Seuss verbatim, raising questions regarding copyright.

Legal landscape. There are several ongoing cases that will shape the legal landscape of LLMs. The New York Times v. Microsoft Corp. lawsuit [[NYT-MC, 2023](#)] represents a pivotal legal challenge to generative AI and copyright, accompanied by similar cases, including Chabon v. OpenAI [[Chabon-OpenAI, 2023](#)] where authors claiming copyright infringement in model training, and Doe v. GitHub [[DOE-GitHub, 2022](#)] where plaintiffs allege verbatim code reproduction in violation of the DMCA. The Times claimed that OpenAI’s model not only trained on its copyrighted articles without proper license but also that ChatGPT could regurgitate substantial passages of those articles on request. In [NYT-MC \[2023\]](#), the publishers argue this behavior amounts to mass copyright infringement since the model’s outputs mimic and even compete with the Times’ content. While researchers such as [Freeman et al. \[2024\]](#) may provide their analysis regarding the veracity of lawsuit [[NYT-MC, 2023](#)]’s copyright infringement claims, such high-profile legal action underscores how memorization is no longer only a technical but a business and public policy problem.

Open Questions

1. Can LLMs be trained to avoid memorizing copyrighted content?
2. How can memorization metrics inform legal analysis (fair use, copyright)?
3. What technical threshold defines legally significant memorization?

8 Conclusion

Memorization in LLMs represents a double-edged sword in the advancement of AI systems. LLMs demonstrate significant capacity to retain and reproduce training data, with larger models exhibiting higher memorization rates of verbatim content, PII, and copyrighted material. This phenomenon creates a fundamental tension between model utility and privacy/legal concerns that must be addressed.

Detection methods for memorization, like prefix-based extraction and membership inference attacks, remain limited, requiring new approaches that work without training data access and standardized auditing tools for comprehensive evaluation. While mitigation strategies, including de-duplication, scrubbing, differential privacy, and model unlearning, show potential, we must better understand their impact on model utility. Ongoing legal challenges highlight the need for balanced frameworks that preserve beneficial knowledge retention while preventing information leakage and establishing appropriate memorization boundaries.

By advancing our understanding in these areas, we can work toward building more powerful and trustworthy AI systems that maintain high performance while respecting privacy and intellectual property concerns. This balance will be essential as LLMs continue to be integrated into increasingly sensitive applications.

References

- Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training, 2018. URL <https://openai.com/research/language-unsupervised>.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman, Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and efficient foundation language models, 2023. URL <https://arxiv.org/abs/2302.13971>.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL <https://arxiv.org/abs/2201.11903>.
- Shuroug A Alowais, Sahar S Alghamdi, Nada Alsuhebany, Tariq Alqahtani, Abdulrahman I Alshaya, Sumaya N Almohareb, Atheer Aldairem, Mohammed Alrashed, Khalid Bin Saleh, Hisham A Badreldin, et al. Revolutionizing healthcare: the role of artificial intelligence in clinical practice. *BMC medical education*, 23(1):689, 2023.
- Diksha Khurana, Aditya Koli, Kiran Khatter, and Sukhdev Singh. Natural language processing: state of the art, current trends and challenges. *Multimedia tools and applications*, 82(3):3713–3744, 2023.
- Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. Large language models: A survey, 2025. URL <https://arxiv.org/abs/2402.06196>.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. A survey on evaluation of large language models. *ACM Trans. Intell. Syst. Technol.*, 15(3), March 2024. ISSN 2157-6904. doi: 10.1145/3641289. URL <https://doi.org/10.1145/3641289>.
- Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large language model (LLMs) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, page 100211, 2024a.

- Badhan Chandra Das, M Hadi Amini, and Yanzhao Wu. Security and privacy challenges of large language models: A survey. *ACM Computing Surveys*, 57(6):1–39, 2025.
- Yogesh K Dwivedi, Nir Kshetri, Laurie Hughes, Emma Louise Slade, Anand Jeyaraj, Arpan Kumar Kar, Abdullah M Baabdullah, Alex Koohang, Vishnupriya Raghavan, Manju Ahuja, et al. Opinion Paper: “so what if ChatGPT wrote it?” multidisciplinary perspectives on opportunities, challenges and implications of generative conversational AI for research, practice and policy. *International journal of information management*, 71:102642, 2023.
- Victoria Smith, Ali Shahin Shamsabadi, Carolyn Ashurst, and Adrian Weller. Identifying and mitigating privacy risks stemming from language models: A survey. *arXiv preprint arXiv:2310.01424*, 2023.
- Sara Abdali, Richard Anarfi, CJ Barberan, and Jia He. Securing large language models: Threats, vulnerabilities and responsible practices. *arXiv preprint arXiv:2403.12503*, 2024.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting training data from large language models, 2021. URL <https://arxiv.org/abs/2012.07805>.
- Frank Cremer, Barry Sheehan, Michael Fortmann, Arash N Kia, Martin Mullins, Finbarr Murphy, and Stefan Materne. Cyber risk and cybersecurity: a systematic review of data availability. *The Geneva papers on risk and insurance. Issues and practice*, 47(3):698, 2022.
- Jan Clusmann, Fiona R Kolbinger, Hannah Sophie Muti, Zunamys I Carrero, Jan-Niklas Eckardt, Narmin Ghaffari Laleh, Chiara Maria Lavinia Löffler, Sophie-Caroline Schwarzkopf, Michaela Unger, Gregory P Veldhuizen, et al. The future landscape of large language models in medicine. *Communications medicine*, 3(1):141, 2023.
- Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. Copyright violations and large language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL <https://openreview.net/forum?id=YokfK5V0oz>.
- Baobao Song, Mengyue Deng, Shiva Raj Pokhrel, Qiujun Lan, Robin Doss, and Gang Li. Digital privacy under attack: Challenges and enablers. *arXiv preprint arXiv:2302.09258*, 2023.
- Joshua Freeman, Chloe Rippe, Edoardo Debenedetti, and Maksym Andriushchenko. Exploring memorization and copyright violation in frontier LLMs: A study of the new york times v. openai 2023 lawsuit, 2024. URL <https://arxiv.org/abs/2412.06370>.
- Shangying Hua, Shuangci Jin, and Shengyi Jiang. The limitations and ethical considerations of ChatGPT. *Data intelligence*, 6(1):201–239, 2024.
- John X Morris, Chawin Sitawarin, Chuan Guo, Narine Kokhlikyan, G Edward Suh, Alexander M Rush, Kamalika Chaudhuri, and Saeed Mahloujifar. How much do language models memorize? *arXiv preprint arXiv:2505.24832*, 2025.
- Keng Siau and Weiyu Wang. Artificial intelligence (AI) ethics: ethics of AI and ethical AI. *Journal of Database Management (JDM)*, 31(2):74–87, 2020.
- Fiona Fui-Hoon Nah, Ruilin Zheng, Jingyuan Cai, Keng Siau, and Langtao Chen. Generative AI and ChatGPT: Applications, challenges, and AI-human collaboration, 2023.

- Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Linhao Yu, Yan Liu, Jiaxuan Li, Bojian Xiong, Deyi Xiong, et al. Evaluating large language models: A comprehensive survey. arXiv preprint arXiv:2310.19736, 2023.
- Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects. Authorea Preprints, 1:1–26, 2023.
- Mert Karabacak and Konstantinos Margetis. Embracing large language models for medical applications: opportunities and challenges. Cureus, 15(5), 2023.
- Bonan Min, Hayley Ross, Elior Sulem, Amir Poursan Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. Recent advances in natural language processing via large pre-trained language models: A survey. ACM Computing Surveys, 56(2):1–40, 2023.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. arXiv preprint arXiv:2303.18223, 1(2), 2023.
- Yutao Zhu, Huaying Yuan, Shuting Wang, Jiongnan Liu, Wenhan Liu, Chenlong Deng, Haonan Chen, Zheng Liu, Zhicheng Dou, and Ji-Rong Wen. Large language models for information retrieval: A survey. arXiv preprint arXiv:2308.07107, 2023.
- Yang Liu, Jiahuan Cao, Chongyu Liu, Kai Ding, and Lianwen Jin. Datasets for large language models: A comprehensive survey. arXiv preprint arXiv:2402.18041, 2024a.
- Devon Myers, Rami Mohawesh, Venkata Ishwarya Chellaboina, Anantha Lakshmi Sathvik, Praveen Venkatesh, Yi-Hui Ho, Hanna Henshaw, Muna Alhawawreh, David Berdik, and Yaser Jararweh. Foundation and large language models: fundamentals, challenges, opportunities, and social impacts. Cluster Computing, 27(1):1–26, 2024.
- Mohaimenul Azam Khan Raiaan, Md Saddam Hossain Mukta, Kaniz Fatema, Nur Mohammad Fahad, Sadman Sakib, Most Marufatul Jannat Mim, Jubaer Ahmad, Mohammed Eunus Ali, and Sami Azam. A review on large language models: Architectures, applications, taxonomies, open issues and challenges. IEEE access, 12:26839–26874, 2024.
- Bernd Carsten Stahl and Damian Eke. The ethics of ChatGPT - exploring the ethical issues of an emerging technology. International Journal of Information Management, 74:102700, 2024.
- Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. Explainability for large language models: A survey. ACM Transactions on Intelligent Systems and Technology, 15(2):1–38, 2024.
- Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In Proceedings of the 2021 ACM conference on fairness, accountability, and transparency, pages 610–623, 2021.
- Xuandong Zhao, Lei Li, and Yu-Xiang Wang. Provably confidential language modelling. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 943–955, 2022a.
- Valentin Hartmann, Anshuman Suri, Vincent Bindschaedler, David Evans, Shruti Tople, and

- Robert West. SoK: Memorization in general-purpose large language models. arXiv preprint arXiv:2310.18362, 2023.
- Ali Satvaty, Suzan Verberne, and Fatih Turkmen. Undesirable memorization in large language models: A survey. arXiv preprint arXiv:2410.02650, 2024.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models, 2023. URL <https://arxiv.org/abs/2202.07646>.
- A Feder Cooper, Katherine Lee, James Grimmelmann, Daphne Ippolito, Christopher Callison-Burch, Christopher A Choquette-Choo, Niloofar Miresghallah, Miles Brundage, David Mimno, Madiha Zahrah Choksi, et al. Report of the 1st workshop on generative ai and law. arXiv preprint arXiv:2311.06477, 2023.
- Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary Lipton, and J Zico Kolter. Rethinking LLM memorization through the lens of adversarial compression. Advances in Neural Information Processing Systems, 37:56244–56267, 2024.
- Nikhil Kandpal, Eric Wallace, and Colin Raffel. Deduplicating training data mitigates privacy risks in language models, 2022. URL <https://arxiv.org/abs/2202.06539>.
- Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models, 2023. URL <https://arxiv.org/abs/2311.17035>.
- Daphne Ippolito, Florian Tramer, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher Choquette Choo, and Nicholas Carlini. Preventing generation of verbatim memorization in language models gives a false sense of privacy. In C. Maria Keet, Hung-Yi Lee, and Sina Zarrieß, editors, Proceedings of the 16th International Natural Language Generation Conference, pages 28–53, Prague, Czechia, September 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.inlg-main.3. URL <https://aclanthology.org/2023.inlg-main.3/>.
- Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. Emergent and predictable memorization in large language models, 2023. URL <https://arxiv.org/abs/2304.11158>.
- Jamie Hayes, Marika Swanberg, Harsh Chaudhari, Itay Yona, Ilya Shumailov, Milad Nasr, Christopher A. Choquette-Choo, Katherine Lee, and A. Feder Cooper. Measuring memorization in language models via probabilistic extraction, 2025. URL <https://arxiv.org/abs/2410.19482>.
- Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. Advances in Neural Information Processing Systems, 33:2881–2891, 2020.
- Chiyuan Zhang, Daphne Ippolito, Katherine Lee, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. Counterfactual memorization in neural language models. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL <https://openreview.net/forum?id=67o9UQgTD0>.
- Aneesh Pappu, Billy Porter, Ilya Shumailov, and Jamie Hayes. Measuring memorization in rlhf for code completion, 2024. URL <https://arxiv.org/abs/2406.11715>.

- Kushal Tirumala, Aram H. Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization without overfitting: Analyzing the training dynamics of large language models, 2022. URL <https://arxiv.org/abs/2205.10770>.
- Qinbin Li, Junyuan Hong, Chulin Xie, Jeffrey Tan, Rachel Xin, Junyi Hou, Xavier Yin, Zhun Wang, Dan Hendrycks, Zhangyang Wang, Bo Li, Bingsheng He, and Dawn Song. LLM-PBE: Assessing data privacy in large language models, 2024. URL <https://arxiv.org/abs/2408.12787>.
- Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2022.
- Hirokazu Kiyomaru, Issa Sugiura, Daisuke Kawahara, and Sadao Kurohashi. A comprehensive analysis of memorization in large language models. In Saad Mahamood, Nguyen Le Minh, and Daphne Ippolito, editors, *Proceedings of the 17th International Natural Language Generation Conference*, pages 584–596, Tokyo, Japan, September 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.inlg-main.45/>.
- Zhepeng Wang, Runxue Bao, Yawen Wu, Jackson Taylor, Cao Xiao, Feng Zheng, Weiwen Jiang, Shangqian Gao, and Yanfu Zhang. Unlocking memorization in large language models with dynamic soft prompting. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9782–9796, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.546. URL <https://aclanthology.org/2024.emnlp-main.546/>.
- Eugene Kharitonov, Marco Baroni, and Dieuwke Hupkes. How BPE affects memorization in transformers. *arXiv preprint arXiv:2110.02782*, 2021.
- Weichen Yu, Tianyu Pang, Qian Liu, Chao Du, Bingyi Kang, Yan Huang, Min Lin, and Shuicheng Yan. Bag of tricks for training data extraction from language models. In *International Conference on Machine Learning*, pages 40306–40320. PMLR, 2023.
- Trishita Tiwari and G. Edward Suh. Sequence-level leakage risk of training data in large language models, 2025. URL <https://arxiv.org/abs/2412.11302>.
- Matthew Jagielski, Om Thakkar, Florian Tramer, Daphne Ippolito, Katherine Lee, Nicholas Carlini, Eric Wallace, Shuang Song, Abhradeep Guha Thakurta, Nicolas Papernot, and Chiyuan Zhang. Measuring forgetting of memorized training examples. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=7bJizxLKrR>.
- Danny Leybzon and Corentin Kervadec. Learning, forgetting, remembering: Insights from tracking LLM memorization during training. In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 43–57, 2024.
- Jing Huang, Diyi Yang, and Christopher Potts. Demystifying verbatim memorization in large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 10711–10732, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.598. URL <https://aclanthology.org/2024.emnlp-main.598/>.
- Fatemehsadat Miresghallah, Archit Uniyal, Tianhao Wang, David Evans, and Taylor Berg-Kirkpatrick. An empirical analysis of memorization in fine-tuned autoregressive language

- models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 1816–1826, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.119. URL <https://aclanthology.org/2022.emnlp-main.119/>.
- Shenglai Zeng, Yaxin Li, Jie Ren, Yiding Liu, Han Xu, Pengfei He, Yue Xing, Shuaiqiang Wang, Jiliang Tang, and Dawei Yin. Exploring memorization in fine-tuned language models, 2024. URL <https://arxiv.org/abs/2310.06714>.
- Xiaoyi Chen, Siyuan Tang, Rui Zhu, Shijun Yan, Lei Jin, Zihao Wang, Liya Su, Zhikun Zhang, XiaoFeng Wang, and Haixu Tang. The janus interface: How fine-tuning in large language models amplifies the privacy risks, 2024. URL <https://arxiv.org/abs/2310.15469>.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. Advances in neural information processing systems, 35: 27730–27744, 2022.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. T”ulu 3: Pushing frontiers in open language model post-training. arXiv preprint arXiv:2411.15124, 2024.
- Xuandong Zhao, Zhewei Kang, Aosong Feng, Sergey Levine, and Dawn Song. Learning to reason without external rewards. arXiv preprint arXiv:2505.19590, 2025.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531, 2015.
- Xuandong Zhao, Zhiguo Yu, Ming Wu, and Lei Li. Compressing sentence representation for semantic retrieval via homomorphic projective distillation. In Findings of the Association for Computational Linguistics: ACL 2022, pages 774–781, Dublin, Ireland, May 2022b. Association for Computational Linguistics. URL <https://aclanthology.org/2022.findings-acl.64>.
- Harsh Chaudhari, Jamie Hayes, Matthew Jagielski, Ilia Shumailov, Milad Nasr, and Alina Oprea. Cascading adversarial bias from injection to distillation in language models, 2025. URL <https://arxiv.org/abs/2505.24842>.
- Aly M Kassem, Omar Mahmoud, Niloofar Miresghallah, Hyunwoo Kim, Yulia Tsvetkov, Yejin Choi, Sherif Saad, and Santu Rana. Alpaca against vicuna: Using LLMs to uncover memorization of LLMs. arXiv preprint arXiv:2403.04801, 2024.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In 2017 IEEE symposium on security and privacy (SP), pages 3–18. IEEE, 2017.
- Michael Duan, Anshuman Suri, Niloofar Miresghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. Do membership inference attacks work on large language models?, 2024. URL <https://arxiv.org/abs/2402.07841>.
- Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learning: Analyzing the connection to overfitting. In 2018 IEEE 31st computer security foundations symposium (CSF), pages 268–282. IEEE, 2018.

- Justus Mattern, Fatemehsadat Mireshghallah, Zhijing Jin, Bernhard Schölkopf, Mrinmaya Sachan, and Taylor Berg-Kirkpatrick. Membership inference attacks against language models via neighbourhood comparison. arXiv preprint arXiv:2305.18462, 2023.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. arXiv preprint arXiv:2310.16789, 2023.
- Jie Zhang, Debeshee Das, Gautam Kamath, and Florian Tramèr. Membership inference attacks cannot prove that a model was trained on your data, 2025. URL <https://arxiv.org/abs/2409.19798>.
- Pratyush Maini and Anshuman Suri. Reassessing emnlp 2024’s best paper: Does divergence-based calibration for membership inference attacks hold up?, 2024. URL <https://www.anshumansuri.com/blog/2024/calibrated-mia/>. Accessed March 5, 2025.
- Weichao Zhang, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. Pretraining data detection for large language models: A divergence-based calibration method. arXiv preprint arXiv:2409.14781, 2024.
- Mustafa Safa Ozdayi, Charith Peris, Jack FitzGerald, Christophe Dupuy, Jimit Majmudar, Haidar Khan, Rahil Parikh, and Rahul Gupta. Controlling the extraction of memorized data from large language models via prompt-tuning, 2023. URL <https://arxiv.org/abs/2305.11759>.
- Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. ProPILE: Probing privacy leakage in large language models. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL <https://openreview.net/forum?id=QkLpGxUboF>.
- Kaixuan Huang, Jiacheng Guo, Zihao Li, Xiang Ji, Jiawei Ge, Wenzhe Li, Yingqing Guo, Tianle Cai, Hui Yuan, Runzhe Wang, Yue Wu, Ming Yin, Shange Tang, Yangsibo Huang, Chi Jin, Xinyun Chen, Chiyuan Zhang, and Mengdi Wang. Math-perturb: Benchmarking llms’ math reasoning abilities against hard perturbations, 2025. URL <https://arxiv.org/abs/2502.06453>.
- Cynthia Dwork. Differential privacy. In International colloquium on automata, languages, and programming, pages 1–12. Springer, 2006.
- Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. Foundations and Trends® in Theoretical Computer Science, 9(3–4):211–407, 2014.
- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In Proceedings of the 2016 ACM SIGSAC conference on computer and communications security, pages 308–318, 2016.
- Xuechen Li, Florian Tramer, Percy Liang, and Tatsunori Hashimoto. Large language models can be strong differentially private learners. arXiv preprint arXiv:2110.05679, 2021.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. ICLR, 1(2):3, 2022.
- Hongbin Liu, Lun Wang, Om Thakkar, Abhradeep Thakurta, and Arun Narayanan. Differentially private parameter-efficient fine-tuning for large asr models. arXiv preprint arXiv:2410.01948, 2024b.
- Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A. Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, and

- Huishuai Zhang. Differentially private fine-tuning of language models, 2022. URL <https://arxiv.org/abs/2110.06500>.
- Da Yu, Huishuai Zhang, Wei Chen, Jian Yin, and Tie-Yan Liu. Large scale private learning via low-rank reparametrization, 2021. URL <https://arxiv.org/abs/2106.09352>.
- Olivia Ma, Jonathan Passerat-Palmbach, and Dmitrii Usynin. Efficient and private: Memorisation under differentially private parameter-efficient fine-tuning in language models. arXiv preprint arXiv:2411.15831, 2024.
- Lynn Chua, Badih Ghazi, Yangsibo Huang, Pritish Kamath, Ravi Kumar, Daogao Liu, Pasin Manurangsi, Amer Sinha, and Chiyuan Zhang. Mind the privacy unit! user-level differential privacy for language model fine-tuning. arXiv preprint arXiv:2406.14322, 2024.
- A Feder Cooper, Christopher A Choquette-Choo, Miranda Bogen, Matthew Jagielski, Katja Filippova, Ken Ziyu Liu, Alexandra Chouldechova, Jamie Hayes, Yangsibo Huang, Niloofar Mireshghallah, et al. Machine unlearning doesn’t do what you think: Lessons for generative ai policy, research, and practice. arXiv preprint arXiv:2412.06966, 2024.
- Ken Ziyu Liu. Machine unlearning in 2024, May 2024. URL <https://ai.stanford.edu/~kzliu/blog/unlearning>.
- Jin Yao, Eli Chien, Minxin Du, Xinyao Niu, Tianhao Wang, Zezhou Cheng, and Xiang Yue. Machine unlearning of pre-trained large language models. arXiv preprint arXiv:2402.15159, 2024b.
- Aditya Golatkar, Alessandro Achille, and Stefano Soatto. 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, 2020.
- Jinghan Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma, and Sijia Liu. Model sparsification can simplify machine unlearning. arXiv preprint arXiv:2304.04934, 1(2):3, 2023.
- Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and Minjoon Seo. Knowledge unlearning for mitigating privacy risks in language models. arXiv preprint arXiv:2210.01504, 2022.
- Sungmin Cha, Sungjun Cho, Dasol Hwang, Honglak Lee, Taesup Moon, and Moontae Lee. Learning to unlearn: Instance-wise unlearning for pre-trained classifiers. In Proceedings of the AAAI conference on artificial intelligence, volume 38, pages 11186–11194, 2024.
- Tong Chen, Faeze Brahman, Jiacheng Liu, Niloofar Mireshghallah, Weijia Shi, Pang Wei Koh, Luke Zettlemoyer, and Hannaneh Hajishirzi. Parapo: Aligning language models to reduce verbatim reproduction of pre-training data, 2025. URL <https://arxiv.org/abs/2504.14452>.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. URL <https://arxiv.org/abs/2305.18290>.
- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J Vazquez, Ulisse Mini, and Monte MacDiarmid. Activation addition: Steering language models without optimization. arXiv e-prints, pages arXiv–2308, 2023.
- Manan Suri, Nishit Anand, and Amisha Bhaskar. Mitigating memorization in LLMs using activation steering. arXiv preprint arXiv:2503.06040, 2025.

- Niklas Stoehr, Mitchell Gordon, Chiyuan Zhang, and Owen Lewis. Localizing paragraph memorization in language models. [arXiv preprint arXiv:2403.19851](#), 2024.
- Ting-Yun Chang, Jesse Thomason, and Robin Jia. Do localization methods actually localize memorized data in LLMs? a tale of two benchmarks. [arXiv preprint arXiv:2311.09060](#), 2023.
- Mor Geva, Avi Caciularu, Kevin Ro Wang, and Yoav Goldberg. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. [arXiv preprint arXiv:2203.14680](#), 2022.
- Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. [arXiv preprint arXiv:2104.08696](#), 2021.
- Zhenhong Zhou, Jiuyang Xiang, Chaomeng Chen, and Sen Su. Quantifying and analyzing entity-level memorization in large language models. In [Proceedings of the AAAI Conference on Artificial Intelligence](#), volume 38, pages 19741–19749, 2024.
- Ashwinee Panda, Christopher A Choquette-Choo, Zhengming Zhang, Yaoqing Yang, and Prateek Mittal. Teach LLMs to phish: Stealing private information from language models. [arXiv preprint arXiv:2403.00871](#), 2024.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. [arXiv preprint arXiv:2101.00027](#), 2020.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open corpus of three trillion tokens for language model pretraining research. [arXiv preprint arXiv:2402.00159](#), 2024.
- Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A Lemley, and Percy Liang. Foundation models and fair use. [Journal of Machine Learning Research](#), 24(400):1–79, 2023.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained transformer language models. [arXiv preprint arXiv:2205.01068](#), 2022.
- NYT-MC. The New York Times v. Microsoft Corporation. Complaint, No. 23-CV-11195, Dec. 2023. State District (S.D.) of New York filed on December 27, 2023.
- Chabon-OpenAI. Chabon v. OpenAI, Inc. Complaint, No. 3:23-cv-04625, Sep. 2023. N.D. Cal. filed on September 8, 2023.
- DOE-GitHub. DOE v. GitHub, Inc. Complaint, No. 4:22-cv-06823, Nov. 2022. N.D. Cal. filed on November 3, 2022.