

---

# Position: Bridge the Gaps between Machine Unlearning and AI Regulation

---

Bill Marino<sup>1,\*</sup>

Meghdad Kurmanji<sup>1,\*</sup>

Nicholas D. Lane<sup>1</sup>

## Abstract

The “right to be forgotten” and the data privacy laws that encode it have motivated machine unlearning since its earliest days. Now, some argue that an inbound wave of artificial intelligence regulations — like the European Union’s Artificial Intelligence Act (AIA) — may offer important new use cases for machine unlearning. However, this position paper argues, this opportunity will only be realized if researchers proactively bridge the (sometimes sizable) gaps between machine unlearning’s state of the art and its potential applications to AI regulation. To demonstrate this point, we use the AIA as our primary case study. Specifically, we deliver a “state of the union” as regards machine unlearning’s current potential (or, in many cases, lack thereof) for aiding compliance with various provisions of the AIA. This starts with a precise cataloging of the potential applications of machine unlearning to AIA compliance. For each, we flag the technical gaps that exist between the potential application and the state of the art of machine unlearning. Finally, we end with a call to action: for machine learning researchers to solve the open technical questions that could unlock machine unlearning’s potential to assist compliance with the AIA — and other AI regulations like it.

## 1 Introduction

Since its inception, Machine Unlearning (MU) has been motivated by the so-called “right to be forgotten” (RTBF) [12], which is encoded in data privacy laws like the European Union’s General Data Protection Regulation (GDPR) [35, Art. 17]. Now, a new wave of AI regulations — including but not limited to the European Union’s Artificial Intelligence Act (AIA) [37] — are working their way through the legislative process [8, 6] or have graduated it and gone into effect [24, 7]. As these AI regulations take shape, researchers have begun to explore whether MU can play a role in supporting compliance with them too [79, 91, 46, 94, 89, 60, 64]. However, recent works call into question whether this motivation really holds water [26].

**In this position paper, we argue that MU’s potential to assist compliance with AI regulation will only be realized if researchers close the technical gaps between MU’s state of the art and this prospective new application.**

We use the AIA as an example to support our argument. This starts with a thorough cataloging of the AIA requirements that MU can hypothetically assist compliance with. For each of these potential use cases, we then scrutinize, from a technical perspective, whether the state of the art (SOTA) of MU can really support the hypothesized application, with a special eye towards auditability (which is especially relevant in the regulation context). In many cases, we identify considerable gaps between the two. Finally, we conclude with a

---

\*Equal contribution. <sup>1</sup> University of Cambridge. Correspondence to: [w1m27@cam.ac.uk](mailto:w1m27@cam.ac.uk).

pointed call for the AI research community to take action and fill these gaps in order to help make MU a more viable tool for assisting AI regulation compliance.

## 2 Machine Unlearning

To set the stage for our analysis, in this section, we define and provide an overview of MU and its key concepts:

### 2.1 Formal Definition of Unlearning

Let  $M = A(D)$  denote a model trained on dataset  $D$  using algorithm  $A$ . An **unlearning query** specifies a **forget-set**  $D_f \subset D$ , with the **retain-set** defined as  $D_r = D \setminus D_f$ . The goal of an unlearning algorithm  $U$  is to remove the influence of  $D_f$  from  $M$ , yielding an unlearned model  $M_u = U(M; D_f, D_r)$ . Depending on the approach,  $U$  may not require access to  $D_r$  [136].

**Definition 2.1.** Following [47],  $U$  is an  $(\epsilon, \delta)$ -**unlearner** if the distribution of  $U(M; D_f, D_r)$  is  $(\epsilon, \delta)$ -close to that of  $A(D_r)$ . Specifically, two distributions  $\mu$  and  $\nu$  are  $(\epsilon, \delta)$ -close if for all measurable events  $B$ :  $\mu(B) \leq e^\epsilon \nu(B) + \delta$  and  $\nu(B) \leq e^\epsilon \mu(B) + \delta$ .

This definition provides a natural taxonomy for MU algorithms. When  $\epsilon = \delta = 0$ ,  $U$  is termed **exact unlearning**; otherwise, it is referred to as **approximate unlearning**.

### 2.2 Informal Definitions

While Def. 2.1 provides a rigorous formulation of MU, researchers commonly use informal interpretations, typically phrased as **removing  $x$  from  $M$** . However, deriving informal definitions directly from Def. 2.1 can be challenging, as the entity to remove may not be explicitly identifiable. For example, in generative models,  $x$  often corresponds to a fact or concept without explicit representation in  $M$  or  $D$ .

Additionally, MU is broadly applied to various methods, but overly general definitions introduce unnecessary complexity, potentially obstructing clear scientific discourse. Therefore, we restrict MU in this paper to ML techniques that explicitly modify the model’s parameter-set (e.g., deletion and retraining, fine-tuning, parameter addition or removal). This scoped definition allows MU to remain a practical tool for applications such as safeguarding and alignment, while methods like guardrailing (or “output suppression” as per Cooper et al. [26]) remain distinct, meriting separate evaluation in regulatory and other contexts.

### 2.3 Evaluation metrics

While Def. 2.1 is widely accepted in the MU community, it presents several challenges in practical settings. First, some works question whether this definition is necessary or sufficient to achieve true MU [111]. Second, in large-scale applications, it is computationally infeasible to directly measure the closeness between the distributions  $A(D_r)$  and  $U(M; D_f, D_r)$ . As a result, researchers often resort to alternative proxies to verify MU. These proxies include performance metrics (e.g., classification accuracy [50] or generative performance using metrics like ROUGE for large language models [88]) and privacy attacks, such as membership inference attacks [61, 113].

### 2.4 Trade-offs and risks

MU algorithms strive balance three key objectives: **Model Utility**, **Forgetting Quality**, and **Efficiency**. In certain privacy-centric applications, forgetting can be synonymous with achieving privacy [83]. Hyperparameters and regularizers impact these trade-offs. For example, in MU via **Fine-tune**, the number of steps and learning rate dictate the balance between forgetting quality and efficiency [127]. Similarly, in **Gradient Ascent**, the number of steps determines the trade-off between effective MU and preserving model’s utility [74].

Additionally, forgetting may sometimes conflict with privacy due to two phenomena. First, unlearning specific data points can inadvertently expose information about others in the

retained set due to the “onion effect” of privacy [14]. Second, over-forgetting [74] a data point may reveal its membership in the original training set—a phenomenon known as the “Streisand Effect” [50]. Addressing these challenges requires careful calibration of MU methods to ensure a delicate equilibrium across these competing objectives.

Beyond these trade-offs, MU introduces risks associated with *untrusted parties* [77] and *malicious unlearning* [100]. Malicious entities could exploit MU to make fake deletion queries, or introduce computation overhead to systems [64].

### 3 The EU’s Artificial Intelligence Act

The AIA sets forth requirements for AI systems and models placed on the market or put into service in the EU [37, Art. 1]. These requirements largely target two categories of AI: AI systems and general-purpose AI (“GPAI”) models. Here, we define these categories and, for each, review some of the AIA requirements that relate to the discussion at hand.

#### 3.1 AI Systems

The AIA broadly defines AI systems to include any “machine-based system that is designed to operate with varying levels of autonomy and that may exhibit adaptiveness after deployment, and that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments” [37, Art. 3.1]. An example of a system that might meet this criteria is ChatGPT [41].

In laying out its rules for these AI systems, the AIA relies on a “risk-based” approach [87], by which an AI system’s perceived degree of risk determines the exact rules that apply to it. Here, the strictest requirements — and the ones most relevant to our discussion — are reserved for those AI systems deemed to be *high-risk* [37, Art. 6]. Such high-risk AI (“HRAI”) systems are subject to a bevy of requirements [37, Chap. III]. Among them, the following are the most relevant to our position:

**Risk management:** HRAI systems must implement risk management systems that identify known and reasonably foreseeable risks that the system may pose to health and safety or to fundamental rights [37, 71, Art. 9.2.a]. Here, risks to *health and safety* includes risks to mental and social well-being as well as physical safety. [2, 33]. Meanwhile, risks to *fundamental rights* includes, among other things, the right to non-discrimination [34], including from biased results [3]. Importantly, wherever these risks are identified, they should be “reasonably mitigated or eliminated through the development or design” of the AI system [37, Art. 9.2-3].

**Accuracy and cybersecurity:** HRAI systems must be designed and developed so as to achieve an “appropriate level” of accuracy and cybersecurity [37, Art. 15.1]. In its Recitals, the AIA stresses that these appropriate levels are a function of the system’s intended purpose as well as the SOTA [37, Rec. 74]. When it comes to cybersecurity, the AIA specifically requires that HRAI systems be “resilient against attempts by unauthorised third parties to alter their use, outputs or performance by exploiting system vulnerabilities” [37, Art. 15.5] and take technical measures to “prevent, detect, respond to, resolve and control for ... data poisoning” as well as “confidentiality attacks” [37, Art. 15.5].

#### 3.2 GPAI models

In contrast to an AI system, a GPAI model is defined as any AI model that is “trained with a large amount of data using self-supervision at scale, that displays significant generality and is capable of competently performing a wide range of distinct tasks regardless of the way the model is placed on the market and that can be integrated into a variety of downstream systems or applications, except AI models that are used for research, development or prototyping activities before they are placed on the market” [37, Art. 3.63]. Some see this as being synonymous with “foundation model” [1]. An example of a GPAI model that might meet this criteria is GPT 3.5, the model powering ChatGPT [41].

In laying out its requirements for GPAI models, the AIA again uses a risk-based approach, with the strictest requirements reserved for GPAI models deemed to carry *systemic risk* [37, Art. 55]. This is defined as the risk of “having a significant impact on the [EU] market due to [its] reach, or due to actual or reasonably foreseeable negative effects on public health, safety, public security, fundamental rights, or the society as a whole, that can be propagated at scale across the value chain” [37, Art. 2.65; Annex III]. This status can be established through proxies, including performance on benchmarks and the amount of compute used during training [37, Art. 51]. While the AIA itself does not further elaborate on what constitutes systemic risk, a companion piece to the AIA posits that it covers risks related to: (1) cyber offense; (2) chemical, biological, radiological, and nuclear (CBRN); (3) loss of control; (4) automated use of models for AI research and development; (5) persuasion and manipulation; and (6) large-scale discrimination [38].

Among the AIA’s requirements for GPAI models that display systemic risk — and those that don’t — the following are the most relevant to our analysis:

**Copyright:** All GPAI model providers must “put in place a policy to comply with Union law on copyright and related rights” [37, Art. 53.c]. Among other things, this policy must respect rightsholders’ requests, per [36, Art. 4.3], to opt out of text and data mining (TDM) on their copyrighted works [37, Rec. 105, Art. 53.c].

**Systemic risk:** GPAI models that display systemic risk must “mitigate” it [37, Art. 55.a-b].

**Cybersecurity:** GPAI models with systemic risk are additionally required to “ensure an adequate level of cybersecurity” [37, Art. 55(d)].

## 4 MU for AIA compliance: a catalog

This Section comprehensively catalogs the potential applications of MU to assist AIA compliance. For each, we analyze the SOTA and its ability to support the potential application, then identify any open questions the research community must resolve in order to bridge the gap between the (including as regards verification, which can be especially important in the regulatory context — for regulators or any other auditors). In sum, we find that the potential applications of MU to assist AIA compliance ultimately roll up into just six separate applications (Fig. 1):

- **Accuracy:** Improve accuracy per EU [37, Arts. 9, 15];
- **Bias:** Mitigate bias per EU [37, Arts. 9, 55];
- **Privacy Attack:** Mitigate confidentiality attacks per EU [37, Arts. 9, 15, 55];
- **Data Poisoning:** Mitigate data poisoning per EU [37, Art. 15];
- **GenAI risk:** Mitigate other risks of generative outputs per EU [37, Arts. 9, 55];
- **Copyright:** Aid compliance with copyright laws, per EU [37, Art. 53].

We study six applications of MU for AIA compliance; in practice, these are often overlapping or interdependent rather than disjoint. For example, in terms of overlap, one unlearning algorithm such as retraining can be used to address multiple use cases: e.g., accuracy, biases, and copyright. In terms of interdependencies, unlearning mislabeled or stale samples to raise overall accuracy (Sec. 4.1) can shift subgroup error profiles, impacting fairness (Sec. 4.2). Differently, debiasing forget-sets may alter susceptibility to membership inference (Sec. 4.3). Likewise, removing backdoored data (Sec. 4.4) often trades off with clean accuracy (Sec. 4.1).

### 4.1 Accuracy

Two AIA provisions may compel HRAI systems providers towards higher levels of accuracy. First, HRAI systems must achieve a level of accuracy appropriate to their intended use and the SOTA [37, Art. 15.1; Rec. 74]. Second, HRAI systems’ risk management practices must include measures to mitigate or eliminate risks to health and safety [37, Art. 9], which, in some domains like medicine, could potentially stem from low accuracy [69, 68, 62]. In either

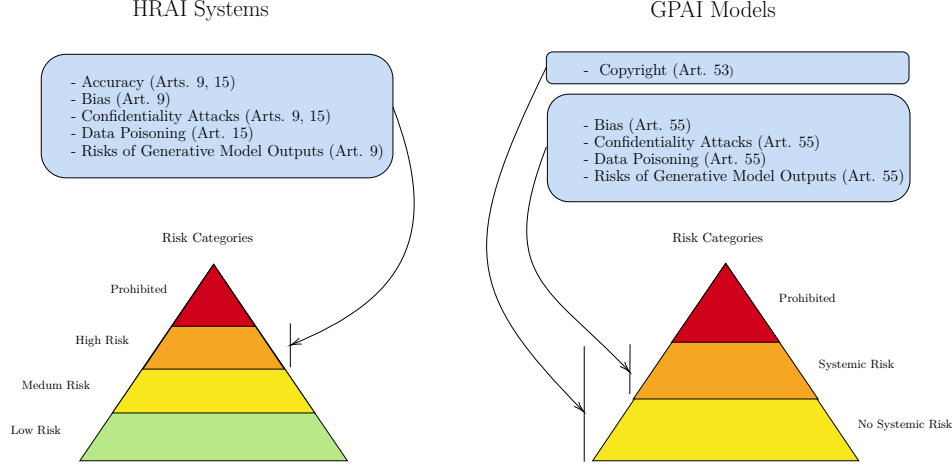


Figure 1: AIA Uses Cases for Machine Unlearning.

case, MU can hypothetically boost accuracy by removing the effect of problematic (e.g., mislabeled) data from the model [110], thus assisting compliance.

This accuracy use case should not require privacy guarantees on the unlearned data [48], because the goal is strictly to boost accuracy to the level deemed appropriate [37, Art. 15] or until the overall residual risk to health and safety posed by the inaccuracy is judged to be "acceptable" [37, Art. 9]. In measuring that, AI providers will presumably account for any inadvertent, counteracting degradation in accuracy caused by the MU itself [76, 11, 125].

**Current SOTA** MU theoretically offer paths towards improving model accuracy by forgetting mislabeled [48, 110, 13, 16, 56], out-of-date and outlier training data points [73, 126, 125, 75], or, potentially, removing noise from medical data [99, 30]. The largest hurdle for this use case might be identifying all of the data points that are leading to inaccuracy (e.g., the mislabeled examples), which can be difficult [49]. It may be good enough to identify only a subset of these examples — so long as accuracy is boosted to levels deemed "appropriate" in light of the intended purpose as well as the SOTA [37, Art. 15.1; Rec. 74]. MU based on subset forget sets have shown success in boosting accuracy; however, other studies have suggested that you need all of polluted data, not just some of this, or it might backfire [49]. It is also important to note that verifying unlearning success is application dependent — and that approximate unlearning should not be expected to yield higher accuracy than exact retraining without the low-quality data.

**Key Points:** (i) Multiple AIA requirements may benefit from MU. (ii) Theoretical guarantees may not be needed. (iii) Evaluation measure is application-dependent.

**Open Problems:** (i) Lack of reliable methods for identifying problematic data to unlearn. (ii) Lack of controllability over trade-offs.

## 4.2 Bias

Providers of both HRAI systems and GPAI models with systemic risk must mitigate certain types of model bias. The former must take measures to mitigate or eliminate risks to fundamental rights, which includes the right to non-discrimination [37, Art. 9]. The latter must take steps to mitigate their models' systemic risk [37, Art. 55], which includes the risk of large-scale discrimination [38]. Bias can occur because unrepresentative or incomplete data prevent the model from perform fairly on different groups or, in the case of generative models, cause it to produce stereotyped or otherwise discriminatory outputs [42]. In all cases, MU can ostensibly help forget the data points or training data patterns causing the bias [98, 75, 104, 72, 16, 138]. An important limiting factor on this use case is that training data that is not there to begin with cannot be forgotten; if the bias is due to a data *deficit*, MU

will not help. Because the goal here is to reduce or eradicate bias, success should ultimately be measured and verified using traditional bias metrics like the difference in performance on various subgroups [27] or, in the case of generative models, the propensity for biased outputs as measured with benchmarks [96].

**Current SOTA** Model debiasing has a longer history than MU [130]. Recently, both exact [58]) and approximate [20, 95, 58]) MU solutions have been offered to mitigate model biases. In the debiasing literature, solutions include *pre-processing*, *in-processing*, and *post-processing* methods [90]. MU, can mainly be considered as a post-processing method. However, it is difficult to draw a separating line between debiasing and MU methods. MU works usually re-use some of the evaluation metrics in the debiasing literature, however, how to evaluate bias is, generally, considered an “open problem” [102]. In order to preserve accuracy (also required under the AIA) by not forgetting data points holistically, [124] use MU to forget only those features that lead to bias.

**Key Points:** (i) MU may aid compliance with multiple bias-related AIA requirements. (ii) MU is only subtractive and never additive, limiting its application to this use case. (iii) De-biasing solutions are not limited to MU.

**Open Problems** (i) Lack of methods for identifying bias counterfactuals. (i) Lack of controllability over trade-offs. (iii) Difficulty of guaranteeing full unlearning of biases, due to generalization.

### 4.3 Confidentiality attacks

The AIA requires providers of both HRAI systems and GPAI models with systemic risk to resolve and control for confidentiality attacks. Providers of HRAI systems must ensure their systems achieve an “appropriate level” of cybersecurity given the intended use and the SOTA, including by taking technical measures to prevent, detect, respond to, resolve and control for confidentiality attacks [37, Art. 15.5; Rec. 74]. Meanwhile, providers of GPAI models with systemic risk must ensure those models reflect “an adequate level of cybersecurity” [37, Art. 55.d], which presumably also includes defending against confidentiality attacks. While the AIA does not define confidentiality attacks, we take them to include any attacks, including data reconstruction and membership inference attacks, that cause a model to reveal confidential details about its training such as data points or membership [115, 23]. This may include confidential training data memorized by generative models [26, 54, 138]. Where such attacks occur — or where there is reason to think they might — MU can ostensibly help defend against them by forgetting the confidential information vulnerable to attack [64, 75, 14, 18, 126, 102, 5]. For this use case, the measure of success (i.e., verification) should be whether confidentiality attacks succeed in the wake of the MU [53], though use case-specific metrics have been developed [88]. When it comes to this use case, there are, importantly, other viable options for protecting against confidentiality attacks, including training with differential privacy (DP) [117, 70].

**Current SOTA.** Multiple MU techniques have been proposed to mitigate confidentiality attacks (or the related problem of inadvertent model leakage of personal data) [28, 4, 85, 17, 119, 12, 11, 10]. As is, applying MU to this use case can carry sizable trade-offs. For example, unlearning some data points for the sake of protecting them from recovery by attackers can jeopardize the privacy of other data points that neighbor the unlearned ones [14] or even increase the risk of membership inference attacks that recover the unlearned data points [19, 5, 74]. Differently, approximate unlearning, when used to delete particular data points, can carry a bias trade-off [131, 95] and an accuracy trade-off that rises as more data is forgotten [52, 88]; either trade-off can potentially undermine AIA compliance. It is also important to note that current MU methods usually fail on new emergent attacks that are devised with new assumptions [135, 66].

**Key Points** (i) MU may aid compliance with multiple confidentiality attack-related AIA requirements. (ii) Due to attack diversity, success should be measured on case-by-case basis. (iii) DP is a strong alternative to MU for this use case.

**Open Problems** (i) Difficulty of providing formal guarantees of attack susceptibility. (ii) Difficulty of applying MU to new, emergent attacks. (iii) Identifying, localizing, and measuring memorization of confidential data.

#### 4.4 Data poisoning

In data poisoning, specially-crafted data points are injected into a training set to alter (e.g., degrade or bias) model behavior to the attacker’s benefit [9]. Backdoor attacks are a type of data poisoning where the injected data points create “triggers” the attacker can exploit during inference [80]. The AIA obligates the providers of both HRAI systems and GPAI model with systemic risk to address such attacks. HRAI system providers must ensure their systems achieve an “appropriate level” of cybersecurity, including via technical measures to “prevent, detect, respond to, resolve and control for” data poisoning attacks [37, Art. 15.5]. Providers of GPAI models with systemic risk, meanwhile, must “ensure an adequate level of cybersecurity” in their models [37, Art. 55.d], which presumably also includes defenses against data poisoning. Where it is known that data poisoning has (or could) occur, MU may help remove the effects of the poisoned data points on the model and, thus, help satisfy these requirements [125, 84, 12, 13, 107]. When it comes to measuring and verifying success for this use case, because the “primary goal is to unlearn the adverse effect due to the manipulated data,” the ideal benchmark would seem to be whether those adverse effects — be they vulnerability to backdoor triggers, bias, or lower accuracy — are eliminated or reduced [49]. For example, Goel et al. [49] measure MU efficacy based on whether proper accuracy on backdoor triggers is restored.

**Current SOTA** Though some work has demonstrated MU can succeed for this use case [120, 105] other works question the effectiveness of using MU to address data poisoning or backdoor attacks specifically [55, 112, 97, 126]). As always, identifying the full forget set (here, the poisoned samples) remains challenging [49]. Some methods, moreover, can have a significant accuracy trade-off on this use case [97]. Such trade-offs can be particularly difficult as poisoned data overlaps with the clean data and, in most cases, they are even visually indistinguishable from each other.

**Key Points** (i) MU may aid compliance with several data poisoning-related AIA requirements. (ii) A proper benchmark should measure the elimination of adverse effects.

**Open Problems** (i) Finding contaminated data at scale is challenging. (ii) Unlearning the backdoor pattern without hurting unaffected data is challenging. (iii) Current MU methods mostly fail on data poisoning use case.

#### 4.5 Other risks of generative outputs

Generative outputs may pose risks to health, safety, and human rights or pose systemic risk that providers of HRAI systems and GPAI models, respectively, must mitigate. For example, HRAI systems’ risk management systems must strive to mitigate or eliminate risks the system poses to health, safety, and fundamental rights [37, Art. 9]. Generative outputs may pose risks to health and safety, e.g., by issuing bad medical advice [122, 59], and may pose risks to the fundamental right of non-discrimination, e.g., by producing stereotyping outputs [92]. For GPAI models with systemic risk, providers of such models must mitigate that risk [37, Art. 55], which could be brought on by generative model outputs that display offensive cyber capabilities, knowledge of CBRN, and more [38, 93]. In all these cases, MU may help mitigate the non-compliant outputs by unlearning the data points or even the concepts in the training set that are causing them [138, 26]. Computationally, it may offer advantages even as compared to other popular alignment techniques like reinforcement learning [128]. Measuring success for this use case should arguably be “context dependent” [128]. That is, the best way to verify the MU’s efficacy is to benchmark the exact behavior that we desire to repair [128]. This could utilize existing benchmarks unrelated to MU [5]. Differently, Li et al. [78] propose a benchmark for measuring MU of CBRN knowledge and approaches

that examine the model parameters for remnants of the unlearned concepts have also been proposed [65].

**Current SOTA** Multiple works use MU to curb undesirable generative model outputs [29, 129, 121, 44]. However, the task is difficult, without agreed-upon best practices [26]. Broad concepts like non-discrimination tend to go beyond individual data points, to latent information which is not easily embodied as a discrete forget set [26, 81]. Even if data points that are intrinsically harmful (e.g., the molecular structure of a bioweapon) are removed, models may still assemble dangerous outputs from latent information in the rest of the dataset [26, 114]. Trying to remove that latent knowledge can risk model utility [26]. As a separate but related issue, AI systems in these scenarios may be dual-use, where the appropriateness of outputs depends on downstream context; this, too, makes identifying the forget set difficult and increases the likelihood of a utility trade-off as the model forgets desirable knowledge alongside undesirable knowledge [26, 108, 102]. All of these issues, in turn, make it difficult if not impossible to specify formal guarantees on the MU [81].

**Key Points** MU may aid compliance with several AIA requirements related to generative outputs.

**Open Problems** (i) Defining the forget set when what we seek to forget is conceptual. (ii) Difficulty of guaranteeing full unlearning of unwanted behaviors, due to generalization. (iii) Mitigating the forgetting of useful knowledge alongside undesirable knowledge.

## 4.6 Copyright

All GPAI model providers must have a policy for complying with EU copyright law [37, Art. 53.c]. Among other things, this policy must honor the TDM opt-outs of rightsholders [37, Art. 53.c; Rec. 105], which is often a feature of AI training [103, 57]. When it comes to AI and copyright law, a distinction is sometimes made between the “input” (training) phase and the “output” (inference) phase of the AI life cycle [103, 101]. At this point in time, the primary compliance risk during the input phase seems to be that an AI training set could include data points that violate TDM opt-outs. When this happens, we assume that using MU to remove the opt-out data points from the trained model does not cure the violation, since the violation occurred at the moment the opt-out data was used for training. That said, MU may still represent a valuable component of a copyright-compliance policy by helping prevent, at the “output” phase, further violations of copyright law when the opt-out data points — or any other copyright-protected data points in the training set — are reproduced to some degree in model outputs [103]. This is a real risk with generative models, which often memorize training data [25, 15]. When MU is applied to this use case, we may measure success by tracking how likely the model is to generate works that are sufficiently similar to the copyrighted works. For example, we might rely on existing benchmarks that measure the tendency of models to produce copyrighted materials [82, 21]. Differently, Ma et al. [86] produce a benchmark for the success of MU in the copyright context.

**Current SOTA** Wu et al. [123] unlearn copyrighted works from diffusion models. At first glance, exact MU would seem to provide a guarantee that copyrighted works in the training set will not be reproduced in outputs [81]. But the fact is that retraining from scratch without the copyrighted data may not be a bulletproof solution for preventing copyright infringement in outputs because substantially similar representations of copyrighted “expressions” (e.g., images of characters like Spiderman) could still appear in outputs based on how the model generalizes from the latent information extracted from the rest of the training set [26]. For the same reason, approximate unlearning aimed at removing the influence of the copyright data points on the model, on top of being hard to prove [81], also cannot ensure that copyrights are not infringed by outputs. In general, the SOTA of approximate unlearning has been deemed “insufficient” for the copyright use case, which may be why practitioners currently lean towards pre- and post-processing tools like prompting and moderation to bring AI into compliance with these laws [81, 109]. Dou et al. [31] “unlearn” copyrighted materials in LLM pre-training datasets by identifying and removing specific weight updates in the model’s parameters that correspond to copyrighted content, evaluating their method by measuring



the similarities between the model’s outputs and the original content. The task of measuring whether substantially similar outputs are being produced is challenging [26].

**Key Points** (i) MU does not help with TDM opt-out violations; the damage is already done. (ii) MU may, however, help with downstream copyright violations in outputs. (iii) To avoid malicious unlearning, TDM opt-outs will have to be verified.

**Open Problems** (i) Difficulty in identifying copyright-infringing works in a dataset. (ii) Difficulty of verifying whether model output owes to copyrighted data or generalization. (iii) Localizing and measuring memorization of copyrighted data is itself an open problem.

## 5 Discussion

MU might offer a path towards compliance with some AIA requirements, but it is not a silver bullet. Throughout this work, we have balanced enthusiasm for MU’s capabilities with a clear-eyed view of its limitations. A recurring challenge across use cases — such as accuracy, bias, and confidentiality — is the difficulty of identifying and isolating harmful or low-quality data. In modern AI models, such information is often encoded in distributed representations, making precise removal difficult and risking forgetting useful knowledge.

In many cases, the target of unlearning (e.g., a fact or concept) lacks a discrete representation. Still, recent work in generative models shows promise: concept editing in diffusion models [63], data attribution [116], and inversion-based techniques [45] all offer ways to trace and remove implicit or emergent representations.

While some applications — like correcting mislabeled data to improve accuracy [73] — are feasible with today’s methods, others (e.g., bias mitigation or copyright control) face steeper barriers. In some cases, MU may be an unnecessarily complex solution relative to alternatives (discussed further in Sec. 6). However, overlaps between applications (e.g., boosting both fairness and accuracy) suggest that well-designed MU interventions could serve multiple regulatory goals simultaneously.

One challenge that stands out is verification or auditability [111]. Throughout the paper we stress that auditability, i.e., the ability of parties, including regulators, to inspect and verify the efficacy of unlearning, is central to making MU viable for these AI regulation compliance use cases. However, as discussed, some of today’s approximate MU methods offer only limited guarantees, complicating auditing. Instead, they rely on empirical proxies (e.g., attack success rates, performance recovery, or distributional similarity) rather than formal proofs. Consequently, these approaches lack strong guarantees that forgetting has been achieved or that residual model behavior no longer depends on the forgotten data. Until addressed, this absence of verifiable guarantees undermines both the auditability and the regulatory utility of MU. Going forward, we advocate for the development of formal forgetting guarantees that can underpin regulator-endorsed standards.

## 6 Alternative Views

The arguments against using MU as a tool for compliance with the AIA or other AI regulation would likely point to its shortcomings, trade-offs, and risks as well as the viable substitutes for MU in these scenarios. Some recent works, for example, broadly question whether MU can really achieve its goals, especially in the generative domain [26, 5, 137, 109]. Other works scrutinize MU’s trade-offs around performance, privacy, security, and cost [125, 14, 133, 72, 137, 43]. These factors could reasonably make alternative methods more appealing for the AI regulation use cases highlighted in this position paper [138, 26]. For example, here are alternatives for each AI regulation use case discussed this paper and their possible advantages as compared to MU:

- **Accuracy:** Full retraining or improving data pipelines may be simpler and more reliable for correcting stale or mislabeled data, while MU may serve best as a fallback when those options are impractical [76, 73].

- **Fairness and bias mitigation:** Pre-, in-, and post-processing bias interventions may offer more holistic and verifiable control than MU, which may play a narrower, subtractive role [95, 51].
- **Confidentiality and privacy:** DP and access control may provide stronger, proactive safeguards against leakage, with MU offering only reactive, case-specific data removal after deployment [67, 50].
- **Security and data poisoning:** Robust or certified training and data-sanitization methods may deliver a more proactive defense, while MU may only serve as a reactive remediation step for excising identified poisoned data or backdoors [22, 97, 39].
- **GenAI risk reduction:** Alignment methods such as Reinforcement Learning from Human Feedback (RLHF) or guardrailing may more effectively constrain model behavior, with MU contributing mainly to remove residual unsafe or sensitive representations [5, 88, 118].
- **Copyright and intellectual property:** Dataset governance, licensing verification, fine-tuning, and output filtering may offer preventive compliance, whereas MU may better function as a post-hoc corrective tool for erasing memorized or stylistically infringing content [134, 132, 26].

These alternative approaches come with their own limitations. For instance, while some may consider DP [67] as a strong alternative to MU, several caveats deserve attention. First, DP mechanisms often struggle to balance tight privacy guarantees with acceptable model utility [106]. This trade-off becomes especially pronounced in high-utility applications. Second, unlike traditional privacy settings where protection is applied uniformly across all data points, MU typically targets a specific subset of data—the so-called “forget set.” In large-scale training corpora that combine individually identifiable data with more publicly available content, applying DP globally may offer overly broad protections that are both inefficient and unnecessary [50]. Third, there are use cases where DP is not sufficient or optimal. For instance, if the objective is to remove a harmful or undesired behavior from a generative model (e.g., misinformation, bias, or offensive content), a DP-trained model may still require explicit MU interventions to mitigate such behaviors.

## 7 Conclusion

There are still sizable challenges that must be cleared before MU will be a viable tool for assisting compliance with the AIA (and, by extension, other AI regulations, which tend to feature recurring principles [40, 32]). To realize MU’s potential for these AI regulation use cases, AI researchers should help solve the open technical problems logged by this position paper. Among other things, this includes work on identifying forget set data points, on resolving the privacy and performance trade-offs of MU, and on resolving the particular challenges posed by generative model outputs. Working collaboratively, we can all help unlock MU’s potential to assist compliance with AI regulation and, by extension, help safeguard the important social values these regulations encode.

## 8 Impact Statement

In essence, this position paper suggests research directions that will help MU evolve into a better tool for assisting AI regulation compliance. Because these AI regulations tend to encode important ethical and societal values around health and safety, non-discrimination, and more, we believe the impact of this paper, in striving to advance AI regulation compliance through MU, will ultimately be the advancement of those important values as well.

## Acknowledgements

*Funding.* We wish to acknowledge the Google Academic Research Awards for their support.

## References

- [1] Ada Lovelace Institute. Foundation models and general purpose AI systems: Understanding impacts and implications. Project report, Ada Lovelace Institute, 2024. URL <https://www.adalovelaceinstitute.org/project/foundation-models-gpai/>.
- [2] A. Armstrong, R. Butler, and K. Gambrell. AI and product safety standards under the EU AI Act. Research paper, Carnegie Endowment for International Peace, March 2024. URL <https://carnegieendowment.org/research/2024/03/ai-and-product-safety-standards-under-the-eu-ai-act>.
- [3] L. Arnold. How the European Union’s AI Act provides insufficient protection against police discrimination. *University of Pennsylvania Carey Law School News*, May 2024. URL <https://www.law.upenn.edu/live/news/16742-how-the-european-unions-ai-act-provides>.
- [4] T. Ashuach, M. Tutek, and Y. Belinkov. Revs: Unlearning sensitive information in language models via rank editing in the vocabulary space, 2024. URL <https://arxiv.org/abs/2406.09325>.
- [5] F. Barez, T. Fu, A. Prabhu, S. Casper, A. Sanyal, A. Bibi, A. O’Gara, R. Kirk, B. Bucknall, T. Fist, L. Ong, P. Torr, K.-Y. Lam, R. Trager, D. Krueger, S. Mindermann, J. Hernandez-Orallo, M. Geva, and Y. Gal. Open problems in machine unlearning for AI safety, 2025. URL <https://arxiv.org/abs/2501.04952>.
- [6] J. Beardwood. The Canadian Artificial Intelligence and Data Act and the EU AI Act: Will sanity prevail as they more closely align? – Part 2 — Changes to both Acts bring them closer together... but not too close. *Computer Law Review International*, 25(5):129–137, 2024. doi: [10.9785/cri-2024-250501](https://doi.org/10.9785/cri-2024-250501). URL <https://doi.org/10.9785/cri-2024-250501>.
- [7] R. Bellan. California’s new AI safety law shows regulation and innovation don’t have to clash. *TechCrunch*, October 2025. URL <https://techcrunch.com/2025/10/05/californias-new-ai-safety-law-shows-regulation-and-innovation-dont-have-to-clash/>.
- [8] L. Belli, Y. Curzi, and W. B. Gaspar. AI regulation in Brazil: Advancements, flows, and need to learn from the data protection experience. *Computer Law & Security Review*, 48: 105767, 2023. ISSN 0267-3649. doi: <https://doi.org/10.1016/j.clsr.2022.105767>. URL <https://www.sciencedirect.com/science/article/pii/S0267364922001108>.
- [9] B. Biggio, B. Nelson, and P. Laskov. Poisoning attacks against support vector machines. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress, 2012. URL <http://icml.cc/2012/papers/880.pdf>.
- [10] J. Borkar. What can we learn from data leakage and unlearning for law? *CoRR*, abs/2307.10476, 2023. doi: [10.48550/ARXIV.2307.10476](https://doi.org/10.48550/ARXIV.2307.10476). URL <https://doi.org/10.48550/arXiv.2307.10476>.
- [11] L. Bourtoutle, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot. Machine unlearning. In *42nd IEEE Symposium on Security and Privacy, SP 2021, San Francisco, CA, USA, 24-27 May 2021*, pages 141–159. IEEE, 2021. doi: [10.1109/SP40001.2021.00019](https://doi.org/10.1109/SP40001.2021.00019). URL <https://doi.org/10.1109/SP40001.2021.00019>.
- [12] Y. Cao and J. Yang. Towards making systems forget with machine unlearning. In *2015 IEEE Symposium on Security and Privacy*, pages 463–480, 2015. doi: [10.1109/SP.2015.35](https://doi.org/10.1109/SP.2015.35).
- [13] Y. Cao, A. F. Yu, A. Aday, E. Stahl, J. Merwine, and J. Yang. Efficient repair of polluted machine learning systems via causal unlearning. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security, ASIACCS ’18*, page 735–747, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450355766. doi: [10.1145/3196494.3196517](https://doi.org/10.1145/3196494.3196517). URL <https://doi.org/10.1145/3196494.3196517>.
- [14] N. Carlini, M. Jagielski, C. Zhang, N. Papernot, A. Terzis, and F. Tramèr. The privacy onion effect: Memorization is relative, 2022. URL <https://arxiv.org/abs/2206.10469>.
- [15] N. Carlini, J. Hayes, M. Nasr, M. Jagielski, V. Sehwag, F. Tramèr, B. Balle, D. Ippolito, and E. Wallace. Extracting training data from diffusion models, 2023. URL <https://arxiv.org/abs/2301.13188>.

- [16] J. Chen and D. Yang. Unlearn what you want to forget: Efficient unlearning for LLMs, 2023. URL <https://arxiv.org/abs/2310.20150>.
- [17] K. Chen, Y. Wang, L. Zhao, C. Jiang, H. Mai, Y. Wu, H. Hong, Y. Shen, J. Mo, L.-L. Huang, J. Peng, X. Wang, and Q. Yang. Private data protection with machine unlearning for next-generation networks. *IEEE Open Journal of the Communications Society*, PP:1–1, 01 2024. doi: 10.1109/OJCOMS.2024.3518503.
- [18] K. Chen, Z. Wang, and B. Mi. Private data protection with machine unlearning in contrastive learning networks. *Mathematics*, 12(24), 2024. ISSN 2227-7390. doi: 10.3390/math12244001. URL <https://www.mdpi.com/2227-7390/12/24/4001>.
- [19] M. Chen, Z. Zhang, T. Wang, M. Backes, M. Humbert, and Y. Zhang. When machine unlearning jeopardizes privacy. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security, CCS '21*, page 896–911. ACM, Nov. 2021. doi: 10.1145/3460120.3484756. URL <http://dx.doi.org/10.1145/3460120.3484756>.
- [20] R. Chen, J. Yang, H. Xiong, J. Bai, T. Hu, J. Hao, Y. Feng, J. T. Zhou, J. Wu, and Z. Liu. Fast model debias with machine unlearning, 2023. URL <https://arxiv.org/abs/2310.12560>.
- [21] T. Chen, A. Asai, N. Mireshghallah, S. Min, J. Grimmelmann, Y. Choi, H. Hajishirzi, L. Zettlemoyer, and P. W. Koh. CopyBench: Measuring literal and non-literal reproduction of copyright-protected text in language model generation, 2024. URL <https://arxiv.org/abs/2407.07087>.
- [22] A. E. Cinà, K. Grosse, A. Demontis, B. Biggio, F. Roli, and M. Pelillo. Machine learning security against data poisoning: Are we there yet? *Computer*, 57(3):26–34, Mar. 2024. ISSN 1558-0814. doi: 10.1109/mc.2023.3299572. URL <http://dx.doi.org/10.1109/MC.2023.3299572>.
- [23] CLTC. Adversarial machine learning, 2024. URL <https://cltc.berkeley.edu/aml/>. CLTC Research Guide.
- [24] Colorado GA. Artificial Intelligence Regulation and Disclosure Act, May 2024. URL <https://leg.colorado.gov/bills/sb24-205>. Senate Bill 24-205.
- [25] A. F. Cooper and J. Grimmelmann. The files are in the computer: On copyright, memorization, and generative AI. *Chicago-Kent Law Review*, 2024. forthcoming.
- [26] A. F. Cooper, C. A. Choquette-Choo, M. Bogen, M. Jagielski, K. Filippova, K. Z. Liu, A. Chouldechova, J. Hayes, Y. Huang, N. Mireshghallah, I. Shumailov, E. Triantafyllou, P. Kairouz, N. Mitchell, P. Liang, D. E. Ho, Y. Choi, S. Koyejo, F. Delgado, J. Grimmelmann, V. Shmatikov, C. D. Sa, S. Barocas, A. Cyphert, M. Lemley, danah boyd, J. W. Vaughan, M. Brundage, D. Bau, S. Neel, A. Z. Jacobs, A. Terzis, H. Wallach, N. Papernot, and K. Lee. Machine unlearning doesn’t do what you think: Lessons for generative AI policy, research, and practice, 2024. URL <https://arxiv.org/abs/2412.06966>.
- [27] D. DeAlcala, I. Serna, A. Morales, J. Fierrez, and J. Ortega-Garcia. Measuring bias in AI models: An statistical approach introducing N-Sigma, 2023. URL <https://arxiv.org/abs/2304.13680>.
- [28] G. Dhingra, S. Sood, Z. M. Wase, A. Bahga, and V. K. Madisetti. Protecting LLMs against privacy attacks while preserving utility. *Journal of Information Security*, 15:448–473, 2024. doi: 10.4236/jis.2024.154026.
- [29] O. Dige, D. Arneja, T. Yau, Q. Zhang, M. Bolandraftar, X. Zhu, and F. Khattak. Can machine unlearning reduce social bias in language models? pages 954–969, 01 2024. doi: 10.18653/v1/2024.emnlp-industry.71.
- [30] N. K. Dinsdale, M. Jenkinson, and A. I. Namburete. Unlearning scanner bias for MRI harmonisation in medical image segmentation. In *Medical Image Understanding and Analysis: 24th Annual Conference, MIUA 2020, Oxford, UK, July 15-17, 2020, Proceedings 24*, pages 15–25. Springer, 2020.
- [31] G. Dou, Z. Liu, Q. Lyu, K. Ding, and E. Wong. Avoiding copyright infringement via large language model unlearning, 2024. URL <https://arxiv.org/abs/2406.10952>.

- [32] N. Díaz-Rodríguez, J. Del Ser, M. Coeckelbergh, M. López de Prado, E. Herrera-Viedma, and F. Herrera. Connecting the dots in trustworthy artificial intelligence: From AI principles, ethics, and key requirements to responsible AI systems and regulation. *Information Fusion*, 99:101896, 2023. ISSN 1566-2535. doi: <https://doi.org/10.1016/j.inffus.2023.101896>. URL <https://www.sciencedirect.com/science/article/pii/S1566253523002129>.
- [33] EC. Questions and answers: Coordinated plan on artificial intelligence 2021 review. Press Release QANDA/21/1683, European Commission, April 2021. URL [https://ec.europa.eu/commission/presscorner/detail/en/qanda\\_21\\_1683](https://ec.europa.eu/commission/presscorner/detail/en/qanda_21_1683).
- [34] EU. Charter of fundamental rights of the European Union, 2000. URL <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:12012P/TXT>. Official Journal of the European Communities, C 364/1.
- [35] EU. General data protection regulation (GDPR), April 2016. URL <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32016R0679>. Official Journal of the European Union, L 119/1.
- [36] EU. Directive (EU) 2019/790 of the European Parliament and of the Council on copyright and related rights in the Digital Single Market and amending Directives 96/9/EC and 2001/29/EC, April 2019. URL <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32019L0790>. Official Journal of the European Union, L 130/92.
- [37] EU. Artificial Intelligence Act, March 2024. URL <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206>. Official Journal of the European Union.
- [38] EU AI Office. First draft of the General-Purpose AI Code of Practice. Policy document, European Union, November 2024. Independent expert draft for stakeholder consultation.
- [39] J. Fan, Q. Yan, M. Li, G. Qu, and Y. Xiao. A survey on data poisoning attacks and defenses. In *2022 7th IEEE International Conference on Data Science in Cyberspace (DSC)*, pages 48–55, 2022. doi: 10.1109/DSC55868.2022.00014.
- [40] S. Feldstein. Evaluating Europe’s push to enact AI regulations: How will this influence global norms? *Democratization*, 31(5):1049–1066, 2024. doi: 10.1080/13510347.2023.2196068. URL <https://doi.org/10.1080/13510347.2023.2196068>.
- [41] D. Fernández-Llorca, E. Gómez, I. Sánchez, et al. An interdisciplinary account of the terminological choices by EU policymakers ahead of the final agreement on the AI Act: AI system, general purpose AI system, foundation model, and generative AI. *Artificial Intelligence and Law*, 2024. doi: 10.1007/s10506-024-09412-y. URL <https://doi.org/10.1007/s10506-024-09412-y>.
- [42] E. Ferrara. Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. *Sci*, 6(1), 2024. ISSN 2413-4155. doi: 10.3390/sci6010003. URL <https://www.mdpi.com/2413-4155/6/1/3>.
- [43] L. Floridi. Machine unlearning: Its nature, scope, and importance for a “delete culture”. *Philosophy & Technology*, 36(42), 2023. doi: 10.1007/s13347-023-00644-5.
- [44] M. Fore, S. Singh, C. Lee, A. Pandey, A. Anastasopoulos, and D. Stamoulis. Unlearning climate misinformation in large language models, 2024. URL <https://arxiv.org/abs/2405.19563>.
- [45] R. Gal, Y. Alaluf, Y. Atzmon, O. Patashnik, A. H. Bermano, G. Chechik, and D. Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022.
- [46] J. Geng, Q. Li, H. Woisetschlaeger, Z. Chen, Y. Wang, P. Nakov, H.-A. Jacobsen, and F. Karray. A comprehensive survey of machine unlearning techniques for large language models, 2025. URL <https://arxiv.org/abs/2503.01854>.
- [47] A. Ginart, M. Guan, G. Valiant, and J. Y. Zou. Making AI forget you: Data deletion in machine learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- [48] S. Goel, A. Prabhu, A. Sanyal, S.-N. Lim, P. Torr, and P. Kumaraguru. Towards adversarial evaluations for inexact machine unlearning, 2023. URL <https://arxiv.org/abs/2201.06640>.
- [49] S. Goel, A. Prabhu, P. Torr, P. Kumaraguru, and A. Sanyal. Corrective machine unlearning, 2024. URL <https://arxiv.org/abs/2402.14015>.

- [50] A. Golatkar, A. Achille, and S. 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.
- [51] S. Goyal, Pooja, A. Kumar, N. Rathod, and A. Verma. Comparative analysis of pre-processing, inprocessing and post-processing methods for bias mitigation: A case study on adult dataset. In *2025 12th International Conference on Computing for Sustainable Global Development (INDIACom)*, pages 1–6, 2025. doi: 10.23919/INDIACom66777.2025.11115514.
- [52] L. Graves, V. Nagisetty, and V. Ganesh. Amnesiac machine learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(13):11516–11524, May 2021. doi: 10.1609/aaai.v35i13.17371. URL <https://ojs.aaai.org/index.php/AAAI/article/view/17371>.
- [53] K. Grimes, C. Abidi, C. Frank, and S. Gallagher. Gone but not forgotten: Improved benchmarks for machine unlearning, 2024. URL <https://arxiv.org/abs/2405.19211>.
- [54] K. Gu, M. R. U. Rashid, N. Sultana, and S. Mehnaz. Second-order information matters: Revisiting machine unlearning for large language models, 2024. URL <https://arxiv.org/abs/2403.10557>.
- [55] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019. doi: 10.1109/ACCESS.2019.2909068.
- [56] E. Gündogdu, A. Unal, and G. Unal. A study regarding machine unlearning on facial attribute data. In *18th IEEE International Conference on Automatic Face and Gesture Recognition, FG 2024, Istanbul, Turkey, May 27-31, 2024*, pages 1–5. IEEE, 2024. doi: 10.1109/FG59268.2024.10581972. URL <https://doi.org/10.1109/FG59268.2024.10581972>.
- [57] Hamburg Regional Court. Hamburg Regional Court, Germany [2024]: Robert Kneschke v. LAION e.V., case no. 310 o 227/23, September 2024. URL <https://www.wipo.int/wipolex/en/judgments/details/2381>. Judgment concerning copyright and text and data mining exceptions under the DSM Directive and German law. Part of the 2024 WIPO Intellectual Property Judges Forum collection.
- [58] L. Han, H. Huang, D. Scheinost, M. Hartley, and M. R. Martínez. Unlearning information bottleneck: Machine unlearning of systematic patterns and biases. *CoRR*, abs/2405.14020, 2024. doi: 10.48550/ARXIV.2405.14020. URL <https://doi.org/10.48550/arXiv.2405.14020>.
- [59] T. Han, S. Nebelung, F. Khader, et al. Medical large language models are susceptible to targeted misinformation attacks. *npj Digital Medicine*, 7:288, 2024. doi: 10.1038/s41746-024-01282-7.
- [60] A. Hatua, T. T. Nguyen, F. Cano, and A. H. Sung. Machine unlearning using forgetting neural networks, 2024. URL <https://arxiv.org/abs/2410.22374>.
- [61] J. Hayes, I. Shumailov, E. Triantafillou, A. Khalifa, and N. Papernot. Inexact unlearning needs more careful evaluations to avoid a false sense of privacy. *arXiv preprint arXiv:2403.01218*, 2024.
- [62] W. D. Heaven. Hundreds of AI tools have been built to catch covid. None of them helped. *MIT Technology Review*, July 2021. URL <https://www.technologyreview.com/2021/07/30/1030329/machine-learning-ai-failed-covid-hospital-diagnosis-pandemic/>.
- [63] A. Hertz, R. Mokady, J. Tenenbaum, K. Aberman, Y. Pritch, and D. Cohen-Or. Prompt-to-prompt image editing with cross attention control. URL <https://arxiv.org/abs/2208.01626>, 1, 2022.
- [64] E. Hine, C. Novelli, M. Taddeo, et al. Supporting trustworthy AI through machine unlearning. *Science and Engineering Ethics*, 30:43, 2024. doi: 10.1007/s11948-024-00500-5. URL <https://doi.org/10.1007/s11948-024-00500-5>.
- [65] Y. Hong, L. Yu, H. Yang, S. Ravfogel, and M. Geva. Intrinsic evaluation of unlearning using parametric knowledge traces, 2024. URL <https://arxiv.org/abs/2406.11614>.
- [66] S. Hu, Y. Fu, S. Wu, and V. Smith. Jogging the memory of unlearned models through targeted relearning attacks. In *ICML 2024 Workshop on Foundation Models in the Wild*, 2024.
- [67] Y. Huang and C. L. Canonne. Tight bounds for machine unlearning via differential privacy. *arXiv preprint arXiv:2309.00886*, 2023.

- [68] C. James, J. Ranson, R. Everson, and D. Llewellyn. Performance of machine learning algorithms for predicting progression to dementia in memory clinic patients. *JAMA Network Open*, 4:e2136553, 12 2021. doi: 10.1001/jamanetworkopen.2021.36553.
- [69] K. R. Jongsma, M. Sand, and M. Milota. Why we should not mistake accuracy of medical AI for efficiency. *NPJ Digital Medicine*, 7(1):57, mar 2024. doi: 10.1038/s41746-024-01047-2.
- [70] G. Kaissis, J. Hayes, A. Ziller, and D. Rueckert. Bounding data reconstruction attacks with the hypothesis testing interpretation of differential privacy, 2023. URL <https://arxiv.org/abs/2307.03928>.
- [71] M. E. Kaminski. Legal fictions in the AI Act. *Boston University Law Review*, 103, November 2023. URL <https://www.bu.edu/bulawreview/files/2023/11/KAMINSKI.pdf>.
- [72] S. Keskpai. Machine unlearning. TechSonar report, European Data Protection Supervisor, January 2024. URL <https://edps.europa.eu/techsonar/machine-unlearning>. TechSonar Series.
- [73] M. Kurmanji, E. Triantafillou, and P. Triantafillou. Machine unlearning in learned databases: An experimental analysis, 2023. URL <https://arxiv.org/abs/2311.17276>.
- [74] M. Kurmanji, P. Triantafillou, J. Hayes, and E. Triantafillou. Towards unbounded machine unlearning, 2023. URL <https://arxiv.org/abs/2302.09880>.
- [75] R. Layne. How to make AI ‘forget’ all the private data it shouldn’t have, February 2024. URL <https://www.library.hbs.edu/working-knowledge/qa-seth-neel-on-machine-unlearning-and-the-right-to-be-forgotten>.
- [76] C. Li, H. Jiang, J. Chen, Y. Zhao, S. Fu, F. Jing, and Y. Guo. An overview of machine unlearning. *High-Confidence Computing*, page 100254, 2024. ISSN 2667-2952. doi: <https://doi.org/10.1016/j.hcc.2024.100254>. URL <https://www.sciencedirect.com/science/article/pii/S2667295224000576>.
- [77] L. Li, X. Ren, H. Yan, X. Liu, and Z. Zhang. Pseudo unlearning via sample swapping with hash. *Information Sciences*, 662:120135, 2024.
- [78] N. Li, A. Pan, A. Gopal, S. Yue, D. Berrios, A. Gatti, J. D. Li, A.-K. Dombrowski, S. Goel, L. Phan, G. Mukobi, N. Helm-Burger, R. Lababidi, L. Justen, A. B. Liu, M. Chen, I. Barrass, O. Zhang, X. Zhu, R. Tamirisa, B. Bharathi, A. Khoja, Z. Zhao, A. Herbert-Voss, C. B. Breuer, S. Marks, O. Patel, A. Zou, M. Mazeika, Z. Wang, P. Oswal, W. Lin, A. A. Hunt, J. Tienken-Harder, K. Y. Shih, K. Talley, J. Guan, R. Kaplan, I. Steneker, D. Campbell, B. Jokubaitis, A. Levinson, J. Wang, W. Qian, K. K. Karmakar, S. Basart, S. Fitz, M. Levine, P. Kumaraguru, U. Tupakula, V. Varadharajan, R. Wang, Y. Shoshitaishvili, J. Ba, K. M. Esvelt, A. Wang, and D. Hendrycks. The WMDP benchmark: Measuring and reducing malicious use with unlearning, 2024. URL <https://arxiv.org/abs/2403.03218>.
- [79] N. Li, C. Zhou, Y. Gao, H. Chen, Z. Zhang, B. Kuang, and A. Fu. Machine unlearning: Taxonomy, metrics, applications, challenges, and prospects. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–21, 2025. doi: 10.1109/TNNLS.2025.3530988.
- [80] J. Lin, L. Dang, M. Rahouti, and K. Xiong. ML attack models: Adversarial attacks and data poisoning attacks, 2021. URL <https://arxiv.org/abs/2112.02797>.
- [81] K. Liu. Machine unlearning: What it is and why it matters, 2023. URL <https://ai.stanford.edu/~kzliu/blog/unlearning>.
- [82] X. Liu, T. Sun, T. Xu, F. Wu, C. Wang, X. Wang, and J. Gao. SHIELD: Evaluation and defense strategies for copyright compliance in LLM text generation, 2024. URL <https://arxiv.org/abs/2406.12975>.
- [83] Z. Liu, G. Dou, E. Chien, C. Zhang, Y. Tian, and Z. Zhu. Breaking the trilemma of privacy, utility, and efficiency via controllable machine unlearning. In *Proceedings of the ACM on Web Conference 2024*, pages 1260–1271, 2024.
- [84] Z. Liu, H. Ye, C. Chen, Y. Zheng, and K.-Y. Lam. Threats, attacks, and defenses in machine unlearning: A survey, 2024. URL <https://arxiv.org/abs/2403.13682>.
- [85] T. Lizzo and L. Heck. Unlearn efficient removal of knowledge in large language models, 2024. URL <https://arxiv.org/abs/2408.04140>.

- [86] R. Ma, Q. Zhou, Y. Jin, D. Zhou, B. Xiao, X. Li, Y. Qu, A. Singh, K. Keutzer, J. Hu, X. Xie, Z. Dong, S. Zhang, and S. Zhou. A dataset and benchmark for copyright infringement unlearning from text-to-image diffusion models, 2024. URL <https://arxiv.org/abs/2403.12052>.
- [87] T. Mahler. Between risk management and proportionality: The risk-based approach in the EU’s Artificial Intelligence Act proposal. *The Swedish Law and Informatics Research Institute*, pages 247–270, Mar. 2022.
- [88] P. Maini, Z. Feng, A. Schwarzschild, Z. C. Lipton, and J. Z. Kolter. TOFU: A task of fictitious unlearning for LLMs, 2024. URL <https://arxiv.org/abs/2401.06121>.
- [89] M. A. Manab. Eternal sunshine of the mechanical mind: The irreconcilability of machine learning and the right to be forgotten, 2024. URL <https://arxiv.org/abs/2403.05592>.
- [90] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan. A survey on bias and fairness in machine learning. *ACM computing surveys (CSUR)*, 54(6):1–35, 2021.
- [91] S. Mercuri, R. Khraishi, R. Okhrati, D. Batra, C. Hamill, T. Ghasempour, and A. Nowlan. An introduction to machine unlearning, 2022. URL <https://arxiv.org/abs/2209.00939>.
- [92] L. Nicoletti and D. Bass. Humans are biased. Generative AI is even worse. *Technology + Equality*, June 2023. URL <https://www.bloomberg.com/graphics/2023-generative-ai-bias/>.
- [93] NIST. NIST trustworthy and responsible AI: Artificial intelligence risk management framework – generative artificial intelligence profile. Technical Report NIST AI 600-1, National Institute of Standards and Technology (NIST), 2024. URL <https://doi.org/10.6028/NIST.AI.600-1>.
- [94] A. Oesterling, U. Bhalla, S. Venkatasubramanian, and H. Lakkaraju. Operationalizing the Blueprint for an AI Bill of Rights: Recommendations for practitioners, researchers, and policy makers, 2024. URL <https://arxiv.org/abs/2407.08689>.
- [95] A. Oesterling, J. Ma, F. P. Calmon, and H. Lakkaraju. Fair machine unlearning: Data removal while mitigating disparities. In S. Dasgupta, S. Mandt, and Y. Li, editors, *International Conference on Artificial Intelligence and Statistics, 2-4 May 2024, Palau de Congressos, Valencia, Spain*, volume 238 of *Proceedings of Machine Learning Research*, pages 3736–3744. PMLR, 2024. URL <https://proceedings.mlr.press/v238/oesterling24a.html>.
- [96] A. Parrish, A. Chen, N. Nangia, V. Padmakumar, J. Phang, J. Thompson, P. M. Htut, and S. R. Bowman. BBQ: A hand-built bias benchmark for question answering, 2022. URL <https://arxiv.org/abs/2110.08193>.
- [97] M. Pawelczyk, J. Z. Di, Y. Lu, G. Kamath, A. Sekhari, and S. Neel. Machine unlearning fails to remove data poisoning attacks, 2024. URL <https://arxiv.org/abs/2406.17216>.
- [98] F. Pedregosa and E. Triantafillou. Announcing the first machine unlearning challenge, June 2023. URL <https://research.google/blog/announcing-the-first-machine-unlearning-challenge/>. Blog post.
- [99] D. Preložnik and Ž. Špiclin. Improving brain MRI segmentation with multi-stage deep domain unlearning. In *International Workshop on PRedictive Intelligence In MEDicine*, pages 99–110. Springer, 2024.
- [100] W. Qian, C. Zhao, W. Le, M. Ma, and M. Huai. Towards understanding and enhancing robustness of deep learning models against malicious unlearning attacks. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 1932–1942, 2023.
- [101] J. P. Quintais. Generative AI, copyright and the AI Act. November 2024. URL <https://ssrn.com/abstract=4912701>. Version 2.
- [102] A. Reuel, B. Bucknall, S. Casper, T. Fist, L. Soder, O. Aarne, L. Hammond, L. Ibrahim, A. Chan, P. Wills, M. Anderljung, B. Garfinkel, L. Heim, A. Trask, G. Mukobi, R. Schaeffer, M. Baker, S. Hooker, I. Solaiman, A. S. Luccioni, N. Rajkumar, N. Moës, J. Ladish, N. Guha, J. Newman, Y. Bengio, T. South, A. Pentland, S. Koyejo, M. J. Kochenderfer, and R. Trager. Open problems in technical AI governance, 2024. URL <https://arxiv.org/abs/2407.14981>.



- [103] E. Rosati. Infringing AI: Liability for AI-generated outputs under international, EU, and UK copyright law. *European Journal of Risk Regulation*, page 1–25, 2024. doi: 10.1017/err.2024.72.
- [104] S. Sai, U. Mittal, V. Chamola, et al. Machine un-learning: An overview of techniques, applications, and future directions. *Cognitive Computation*, 16:482–506, 2024. doi: 10.1007/s12559-023-10219-3. URL <https://doi.org/10.1007/s12559-023-10219-3>.
- [105] S. Schoepf, J. Foster, and A. Brintrup. Potion: Towards poison unlearning. *arXiv preprint arXiv:2406.09173*, 2024.
- [106] J. Seeman and D. Susser. Between privacy and utility: On differential privacy in theory and practice. *ACM Journal on Responsible Computing*, 1(1):1–18, 2024.
- [107] S. Shan, A. N. Bhagoji, H. Zheng, and B. Y. Zhao. Poison forensics: Traceback of data poisoning attacks in neural networks. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 3575–3592, Boston, MA, Aug. 2022. USENIX Association. ISBN 978-1-939133-31-1. URL <https://www.usenix.org/conference/usenixsecurity22/presentation/shan>.
- [108] W. Shi, J. Lee, Y. Huang, S. Malladi, J. Zhao, A. Holtzman, D. Liu, L. Zettlemoyer, N. A. Smith, and C. Zhang. MUSE: Machine unlearning six-way evaluation for language models, 2024. URL <https://arxiv.org/abs/2407.06460>.
- [109] I. Shumailov, J. Hayes, E. Triantafillou, G. Ortiz-Jimenez, N. Papernot, M. Jagielski, I. Yona, H. Howard, and E. Bagdasaryan. Ununlearning: Unlearning is not sufficient for content regulation in advanced generative AI, 2024. URL <https://arxiv.org/abs/2407.00106>.
- [110] I. Sugiura, S. Okamura, and N. Yanai. Removing mislabeled data from trained models via machine unlearning. *IEICE Transactions on Information and Systems*, advpub:2024DAT0002, 2024. doi: 10.1587/transinf.2024DAT0002.
- [111] A. Thudi, H. Jia, I. Shumailov, and N. Papernot. On the necessity of auditable algorithmic definitions for machine unlearning. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 4007–4022, 2022.
- [112] V. Tolpegin, S. Truex, M. E. Gursoy, and L. Liu. Data poisoning attacks against federated learning systems. In L. Chen, N. Li, K. Liang, and S. Schneider, editors, *Computer Security – ESORICS 2020*, pages 480–501, Cham, 2020. Springer International Publishing. ISBN 978-3-030-58951-6.
- [113] E. Triantafillou, P. Kairouz, F. Pedregosa, J. Hayes, M. Kurmanji, K. Zhao, V. Dumoulin, J. J. Junior, I. Mitliagkas, J. Wan, et al. Are we making progress in unlearning? findings from the first NeurIPS unlearning competition. *arXiv preprint arXiv:2406.09073*, 2024.
- [114] UK DSIT. International AI safety report: The international scientific report on the safety of advanced AI. Technical report, UK Government, January 2025. URL [https://assets.publishing.service.gov.uk/media/679a0c48a77d250007d313ee/International\\_AI\\_Safety\\_Report\\_2025\\_accessible\\_f.pdf](https://assets.publishing.service.gov.uk/media/679a0c48a77d250007d313ee/International_AI_Safety_Report_2025_accessible_f.pdf).
- [115] A. Vassilev, A. Oprea, A. Fordyce, and H. Anderson. Adversarial machine learning: A taxonomy and terminology of attacks and mitigations. NIST AI 100-2e2023, National Institute of Standards and Technology, 2023. URL <https://nvlpubs.nist.gov/nistpubs/ai/NIST.AI.100-2e2023.pdf>.
- [116] S.-Y. Wang, A. Hertzmann, A. Efros, J.-Y. Zhu, and R. Zhang. Data attribution for text-to-image models by unlearning synthesized images. *Advances in Neural Information Processing Systems*, 37:4235–4266, 2024.
- [117] Y. Wang, Q. Wang, L. Zhao, and C. Wang. Differential privacy in deep learning: Privacy and beyond. *Future Generation Computer Systems*, 148:408–424, 2023. ISSN 0167-739X. doi: <https://doi.org/10.1016/j.future.2023.06.010>. URL <https://www.sciencedirect.com/science/article/pii/S0167739X23002315>.
- [118] Z. Wang, B. Bi, S. K. Pentyla, K. Ramnath, S. Chaudhuri, S. Mehrotra, Zixu, Zhu, X.-B. Mao, S. Asur, Na, and Cheng. A comprehensive survey of LLM alignment techniques: RLHF, RLAI, PPO, DPO and more, 2024. URL <https://arxiv.org/abs/2407.16216>.
- [119] Z. Wang, S. Chen, C. Li, L. Zhao, and Y. Liu. Applying machine unlearning techniques to mitigate privacy leakage in large language models: An empirical study. Sept. 2024. doi: 10.22541/au.172712647.70020033/v1. URL <http://dx.doi.org/10.22541/au.172712647.70020033/v1>.

- [120] A. Warnecke, L. Pirch, C. Wressnegger, and K. Rieck. Machine unlearning of features and labels, 2023. URL <https://arxiv.org/abs/2108.11577>.
- [121] X. Wei, N. Kumar, and H. Zhang. Addressing bias in generative AI: Challenges and research opportunities in information management. *Information & Management*, page 104103, 2025. ISSN 0378-7206. doi: <https://doi.org/10.1016/j.im.2025.104103>. URL <https://www.sciencedirect.com/science/article/pii/S0378720625000060>.
- [122] K. Wu, E. Wu, D. E. Ho, and J. Zou. Generating medical errors: GenAI and erroneous medical references. February 2024.
- [123] Y. Wu, S. Zhou, M. Yang, L. Wang, H. Chang, W. Zhu, X. Hu, X. Zhou, and X. Yang. Unlearning concepts in diffusion model via concept domain correction and concept preserving gradient, 2024. URL <https://arxiv.org/abs/2405.15304>.
- [124] H. Xu, T. Zhu, W. Zhou, and W. Zhao. Don’t forget too much: Towards machine unlearning on feature level, 2024. URL <https://arxiv.org/abs/2406.10951>.
- [125] J. Xu, Z. Wu, C. Wang, and X. Jia. Machine unlearning: Solutions and challenges. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 8(3):2150–2168, June 2024. ISSN 2471-285X. doi: 10.1109/tetci.2024.3379240. URL <http://dx.doi.org/10.1109/TETCI.2024.3379240>.
- [126] Y. Xu. Machine unlearning for traditional models and large language models: A short survey, 2024. URL <https://arxiv.org/abs/2404.01206>.
- [127] J. Yao, E. Chien, M. Du, X. Niu, T. Wang, Z. Cheng, and X. Yue. Machine unlearning of pre-trained large language models. *arXiv preprint arXiv:2402.15159*, 2024.
- [128] Y. Yao, X. Xu, and Y. Liu. Large language model unlearning, 2024. URL <https://arxiv.org/abs/2310.10683>.
- [129] C. Yu, S. Jeoung, A. Kasi, P. Yu, and H. Ji. Unlearning bias in language models by partitioning gradients. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 6032–6048, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.375. URL <https://aclanthology.org/2023.findings-acl.375>.
- [130] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork. Learning fair representations. In *International conference on machine learning*, pages 325–333. PMLR, 2013.
- [131] D. Zhang, S. Pan, T. Hoang, Z. Xing, M. Staples, X. Xu, L. Yao, Q. Lu, and L. Zhu. To be forgotten or to be fair: Unveiling fairness implications of machine unlearning methods. *arXiv preprint arXiv:2302.03350*, 2023.
- [132] D. Zhang, Z. Xu, and W. Zhao. LLMs and copyright risks: Benchmarks and mitigation approaches. In M. Lomeli, S. Swayamdipta, and R. Zhang, editors, *Proceedings of the 2025 Annual Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 5: Tutorial Abstracts)*, pages 44–50, Albuquerque, New Mexico, May 2025. Association for Computational Linguistics. ISBN 979-8-89176-193-3. doi: 10.18653/v1/2025.naacl-tutorial.7. URL <https://aclanthology.org/2025.naacl-tutorial.7/>.
- [133] H. Zhang, T. Nakamura, T. Isohara, and K. Sakurai. A review on machine unlearning. *SN Computer Science*, 4(4), Apr. 2023. ISSN 2661-8907. doi: 10.1007/s42979-023-01767-4. URL <http://dx.doi.org/10.1007/s42979-023-01767-4>.
- [134] J. Zhang, J. Yu, M. Marone, B. V. Durme, and D. Khashabi. Certified mitigation of worst-case LLM copyright infringement, 2025. URL <https://arxiv.org/abs/2504.16046>.
- [135] Z. Zhang, F. Wang, X. Li, Z. Wu, X. Tang, H. Liu, Q. He, W. Yin, and S. Wang. Does your LLM truly unlearn? an embarrassingly simple approach to recover unlearned knowledge. *arXiv preprint arXiv:2410.16454*, 2024.
- [136] K. Zhao, M. Kurmanji, G.-O. Bărbulescu, E. Triantafillou, and P. Triantafillou. What makes unlearning hard and what to do about it. *arXiv preprint arXiv:2406.01257*, 2024.
- [137] S. Zhou, L. Wang, J. Ye, Y. Wu, and H. Chang. On the limitations and prospects of machine unlearning for generative AI, 2024. URL <https://arxiv.org/abs/2408.00376>.
- [138] J. Lucki, B. Wei, Y. Huang, P. Henderson, F. Tramèr, and J. Rando. An adversarial perspective on machine unlearning for AI safety, 2024. URL <https://arxiv.org/abs/2409.18025>.