
Position: Adversarial ML for LLMs Is Not Making Any Progress

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In the past decade, considerable research effort has been devoted to securing
2 machine learning (ML) models that operate in adversarial settings. Yet, progress
3 has been slow even for simple “toy” problems (e.g., robustness to small adversarial
4 perturbations) and is often hindered by non-rigorous evaluations. Today, adversarial
5 ML research has shifted towards studying larger, general-purpose language models.
6 In this position paper, we argue that the situation is now even worse: **in the era of**
7 **LLMs, the field of adversarial ML studies problems that are (1) less clearly**
8 **defined, (2) harder to solve, and (3) even more challenging to evaluate.** As
9 a result, we caution that yet another decade of work on adversarial ML may be
10 failing to produce meaningful progress.

1 1 Introduction

12 When adversarial machine learning emerged as a field, it focused on attacking and defending simple
13 models with well-defined objectives. For example, misclassifying a spam message as safe (Graham-
14 Cumming, 2004) or images in deep learning models (Biggio et al., 2013; Szegedy, 2013; Goodfellow
15 et al., 2014). These early problems were well-defined: the attack goals were clear (e.g., cause a
16 misclassification), the target models were relatively simple (e.g., linear classifiers, small neural
17 networks), the threat models were simple (e.g., perturb pixels by at most 8/255), and the evaluation
18 metrics were straightforward (e.g., accuracy on a test set). Yet the field has struggled to develop
19 robust solutions or even to fully understand why these vulnerabilities exist (Barreno et al., 2006;
20 Shafahi et al., 2019). Even fundamental “toy” problems like robustness to ℓ_p -bounded perturbations,
21 remain largely unsolved to this day, and many defense evaluations still lack rigor (Carlini & Wagner,
22 2017; Carlini et al., 2019; Tramer et al., 2020).

23 Recently, the focus of the field has since shifted towards studying adversarial problems with large
24 language models (LLMs) and other generative models. **In this position paper, we argue that**
25 **these new problems are significantly harder to define, solve and evaluate; making progress**
26 **increasingly difficult to track.**

27 Due to their general-purpose nature, LLMs are not designed to solve any single well-defined “task” to
28 be secured. Instead, the field now considers a more holistic notion of “safety”, with adversarial objec-
29 tives that are hard to define formally (e.g., making an LLM produce “harmful” responses) (Christiano
30 et al., 2017; Ouyang et al., 2022; Bai et al., 2022; Casper et al., 2023). These safety properties are
31 also often considered for unbounded threat models, thereby leading to much stronger adversaries
32 (e.g., with the ability to adversarially fine-tune a model or to prompt it in arbitrary ways). Due to this
33 large attack space—and the difficulty of directly optimizing over it (Carlini et al., 2024)—attacks
34 are increasingly ad-hoc and human driven (Li et al., 2024a). This further complicates the task for
35 defenders, who cannot automatically search over strong, adaptive attacks.

Table 1: Challenges in different research areas when defining and solving adversarial ML problems.

Research Area	Challenges							
	Defining			Solving		Evaluating		
	(\$2.1.1) Defining Success	(\$2.1.2) Bounding Attacks	(\$2.1.3) Delimiting Data	(\$2.2.1) Attack Search	(\$2.2.2) Principled Defenses	(\$2.3.1) Measuring Harm & Utility	(\$2.3.2) Ensuring Reproduc.	
(\$3.1) Jailbreaks	✓	✓		✓	✓	✓	✓	
(\$3.2) Un-finetunable Models	✓	✓		✓	✓	✓	✓	
(\$3.3) Poisoning + Backdoors	✓	✓	✓	✓	✓	✓	✓	
(\$3.4) Prompt Injections	✓	✓		✓	✓	✓	✓	
(\$3.5) Membership Inference	✓		✓					✓
(\$3.6) Unlearning	✓	✓	✓	✓	✓	✓		

36 Beyond making the technical problems harder, we argue that generative models have also made
 37 evaluation and benchmarking of attacks and defenses more challenging. Measuring attack success is
 38 no longer as straightforward as measuring misclassification rates; it instead requires careful (human)
 39 evaluation of possible harms present in natural language outputs (Mazeika et al., 2024; Chao et al.,
 40 2024). In a similar vein, evaluating whether defenses preserve the utility of the original model has
 41 become more nuanced: instead of measuring test accuracy on a single task, we now have to determine
 42 whether a model maintains its general-purpose capabilities (Cui et al., 2024; Mai et al., 2025).

43 Finally, reproducible benchmarking became harder as many state-of-the-art models are deployed
 44 via black-box APIs that may receive constant updates and patches as newer attacks are released. As
 45 these changes are often not reported, reproducing results or making meaningful comparisons between
 46 different approaches becomes nearly impossible.

47 In this position paper, we use several case studies of research areas in adversarial ML to illustrate
 48 the increasing complexity in both attacks and defenses. We first analyze how traditional research
 49 problems have evolved to become fundamentally harder to formally define and solve (Section 2).
 50 We then present case studies that illustrate these new challenges (Section 3). Finally, we discuss
 51 our perspective on why these changes represent a fundamental challenge to progress in the field and
 52 alternative views on the evolution of adversarial ML (Section 4).

53 2 New Challenges in Defining, Solving, and Evaluating Adversarial ML 54 Problems

55 Traditional ML models were designed and trained for specific and narrow tasks—often classification.
 56 For example, computer vision models used to classify images into a fixed set of classes (Krizhevsky
 57 et al., 2012), and natural language processing models used to perform textual analysis on individual
 58 sentences (Richardson et al., 2013; Rajpurkar, 2016). Additionally, the training and test data were
 59 clearly delineated as inputs were discrete and bounded units (individual images or sentences). In
 60 these settings, adversarial objectives could be clearly specified. For example, misclassifying as many
 61 inputs as possible (i.e., adversarial examples (Szegedy, 2013; Goodfellow et al., 2014)) or inferring if
 62 a given data point was used for training (i.e., membership inference (Shokri et al., 2017)).

63 However, LLMs have fundamentally changed this landscape. Models no longer perform narrow tasks
 64 but serve as general-purpose systems that produce free-form and unbounded outputs. As a result,
 65 defining “security” or “safety” properties of the AI system has become more challenging, with the field
 66 focusing on general definitions (e.g., a model should not produce outputs that can “harm others”¹).
 67 Adversarial objectives related to training data (e.g., membership inference or unlearning) have also
 68 become more ill-defined, as the training set(s) of LLMs span virtually the entire Internet (Gao et al.,
 69 2020), with no clear boundaries between data points or between train and test sets.

70 In this section, we identify three core challenges, each split into several sub-challenges, that make
 71 adversarial ML for LLMs *harder to define, harder to solve, and harder to evaluate*. We provide a
 72 summary of the challenges faced in different research areas in Table 1.

¹<https://openai.com/policies/usage-policies/>

73 In Section 3, we elaborate on how these challenges hinder progress by analyzing specific case studies:
74 *Jailbreaks* (Section 3.1), *Un-finetunable Models* (Section 3.2), *Poisoning and Backdoors* (Section 3.3),
75 *Prompt Injections* (Section 3.4), *Membership Inference* (Section 3.5), and *Unlearning* (Section 3.6).

76 2.1 Problems are Harder to Define

77 2.1.1 Defining Success of Attacks and Defenses

78 In the past, adversarial problems for classification models typically involved concrete objectives (e.g.,
79 misclassifying images), which could be easily measured by accuracy on a set of clean or perturbed
80 inputs. Now, the lack of a single well-defined task makes it unclear what criteria constitute a genuine
81 success or failure for attacks or defenses.

82 LLMs produce free-form text in which goals become subjective. Developers now aim to optimize
83 abstract properties like helpfulness, honesty, and harmlessness (Bai et al., 2022), while adversaries
84 may try to obtain generically harmful outputs. Thus, measuring attack success—i.e., whether an
85 output is actually harmful or violates the developer policies—also becomes subjective.

86 2.1.2 Defining and Bounding the Attack Space

87 In prior robustness settings (e.g., with classification models), the adversary was often constrained to
88 perturb inputs within an ℓ_p -ball around a given image. This served as a meaningful *necessary* but *not sufficient*
89 condition for robustness Gilmer et al. (2018), allowing quantitative comparisons of different
90 methods (Goodfellow et al., 2014).

91 For LLMs, researchers almost always allow the search space for attacks to be unbounded, since
92 any input could potentially elicit a violation of a safety property (Wei et al., 2024a). The shift from
93 input-dependent to input-*independent* constraints makes it harder to specify adversarial capabilities
94 that allow us to compare attacks and defenses. Beyond unbounded inputs, threat models have also
95 become more permissive. In traditional adversarial ML problems (e.g., adversarial examples or
96 poisoning), the strongest adversaries had white-box access to model weights, but could not alter
97 the model’s functionality. Now attackers need not maintain the model’s general capabilities as long
98 as they can elicit the desired harmful information, enabling stronger attacks such as fine-tuning or
99 pruning (Qi et al., 2024b; Wei et al., 2024b)².

100 Moreover, the set of attacks that should be ruled out may not always be obvious. While one could say
101 “any input that leads to harmful content is a valid attack,” trivial attacks such as prompting “please
102 repeat [harmful text]” do not reveal meaningful new vulnerabilities. Hence, there is no clear universal
103 standard on what sorts of prompts or transformations count as “valid” or “novel” adversarial inputs.

104 2.1.3 Delimiting Data

105 In many research areas traditionally studied in adversarial ML, such as unlearning or privacy protection,
106 the notion of a *training data point* plays a crucial role. Previously, a model was trained on a
107 carefully curated dataset with strict train/test splits; each data point (such as a single labeled image)
108 was distinct, and known to researchers. In contrast, generative models are trained on vast corpora,
109 where similar, or even identical, content may appear across multiple subsets of the training set. The
110 exact contents of the training data are also rarely publicly released (Nasr et al., 2025). The notion of
111 a held-out (IID) test set no longer really exists.

112 2.2 Problems are Harder to Solve

113 2.2.1 Searching over Attacks

114 The optimization landscape for most adversarial ML problems has become significantly more
115 complex with LLMs. In traditional classification problems, such as crafting adversarial images,
116 the objective function was clear: maximize the loss on the correct prediction while minimizing
117 perturbation size. This objective could be formalized and optimized by propagating gradients to the

²For adversarial robustness in image classifiers, the ability to finetune the victim model would be a trivial attack vector, since the attacker could simply fine-tune the model to have low accuracy.

118 input space (Madry, 2017). These automated attacks outperformed humans and consistently found
119 worst-case attacks (Carlini et al., 2017).

120 However, the attack surface for LLMs is much larger and harder to define (see Section 2.1.2). There
121 is no longer a single well-defined “task”, and safety properties cannot be expressed with formal loss
122 functions—they are qualitative, context-dependent, and often subjective (Bai et al., 2022).

123 Even if we define a “toy” attack objective (e.g., making the model output an affirmative response such
124 as “Sure, I can help you with that” (Zou et al., 2023)), finding good attacks remains hard (Carlini et al.,
125 2024). Discrete text inputs makes gradient-based methods less effective (Carlini et al., 2024; Rando
126 et al., 2024), and the vast search space makes exploration impractical. Perhaps most telling, manual
127 attacks still outperform automated methods at finding worst-case inputs (Li et al., 2024a). Many
128 successful attacks on LLMs exploit qualitative properties that are hard to optimize automatically,
129 such as persona modulation (Shah et al., 2023), multi-turn conversations (Anil et al., 2024), and
130 social engineering techniques (Zeng et al., 2024). In contrast, current optimization methods typically
131 generate gibberish inputs (Zou et al., 2023; Thompson & Sklar, 2024).

132 **2.2.2 Building Principled Defenses**

133 In traditional adversarial tasks, researchers could devise *certified* defenses (Cohen et al., 2019) or well-
134 motivated empirical defenses such as adversarial training (Madry, 2017), where key properties of the
135 problem (like bounded input perturbations) were explicitly understood. Moreover, the performance
136 of these defenses could be evaluated with strong, adaptive white-box attacks (Tramer et al., 2020).

137 In contrast, for LLMs the adversarial objectives are typically not formally defined (see Section 2.1.1)
138 and the attack space is challenging to bound (see Section 2.1.2). As a result, there is little hope to
139 build defenses upon principled foundations. Existing defenses rely on ad-hoc approaches, through
140 either: (1) adversarial training against *known* successful attacks Bai et al. (2022); Wallace et al.
141 (2024); (2) “virtual” adversarial training in the model’s latent space Miyato et al. (2018); Casper et al.
142 (2024b); Sheshadri et al. (2024); (3) building external classifiers or detectors (Inan et al., 2023); (4)
143 or random preprocessing (Robey et al., 2023). Crucially, none of these approaches produce systems
144 whose security can be analyzed or quantified in a well-defined formal. It is thus not too surprising
145 that the original evaluations of some of these defenses overestimate their robustness (Chi et al., 2024;
146 Qi et al., 2024a; Lucki et al., 2024).

147 **2.3 Problems are Harder to Evaluate**

148 **2.3.1 Measuring Attack Harm and Defense Utility**

149 Since safety properties for LLMs are hard to formally define, it has become customary to use LLMs
150 themselves as a fuzzy “judge” to determine harmfulness (e.g., when evaluating jailbreaks or prompt
151 injections (Mazeika et al., 2024)). But this approach suffers from a number of issues. First, judges
152 fall short of human judgment.³ For instance, many implementations often default to considering any
153 non-refusal response as a successful attack even if the content is harmless (Souly et al., 2024). Second,
154 judges themselves may be vulnerable to attacks (Mangaokar et al., 2024; Raina et al., 2024). Third,
155 using LLMs-as-judges to evaluate defenses can create artificial correlations that bias evaluations. For
156 example, a defense that implements an output filter similar to the judge may achieve near-perfect
157 scores without necessarily being effective against prompts where the judge fails (Liu et al., 2024).

158 Measuring benign utility of defenses—whether they preserve other capabilities—is also non-trivial.
159 Unlike classification tasks where accuracy on a fixed test set is standard, LLMs can be used for an
160 open-ended array of tasks. A defense can trivially produce a safe-but-useless model by refusing all
161 requests. Thus, any evaluation framework must somehow account for the model’s usefulness to the
162 end-user, which is subjective and context-dependent (Cui et al., 2024).

³Even (non-expert) humans have a hard time judging harmfulness of model responses, e.g., when judging whether “instructions for building a bomb” truly yield a useful design.

163 **2.3.2 Reproducing and Comparing Results**

164 In earlier, more controlled research environments, practitioners had detailed information about a
165 model’s architecture, training data, and training pipeline, enabling precise definitions of threats,
166 defenses, and success criteria. This transparency made it straightforward to track progress.

167 Many influential LLMs are now closed-source and updated silently over time (Chao et al., 2024),
168 making it unclear which version of a system is being tested. Moreover, instead of investigating
169 a single, well-defined model, one must analyze an entire system that may incorporate multiple
170 pre-processing, post-processing, or other defense mechanisms.

171 This lack of transparency severely undermines reproducibility. Researchers cannot confirm whether
172 observed behaviors persist across different snapshots of the system, nor can they reliably benchmark
173 potential solutions. Consequently, adversarial ML problems become harder to define—let alone
174 solve and evaluate. While black-box or discrete optimization approaches can help reveal some
175 vulnerabilities, they provide only limited insight into the model’s internals, leaving many critical
176 security and privacy questions unanswered (Casper et al., 2024a; Carlini et al., 2024).

177 **3 Case Studies**

178 **3.1 Jailbreaks**

179 Jailbreaks illustrate many of the new challenges in adversarial research. Jailbreaks are adversarial text
180 inputs for language models that bypass safeguards to generate “harmful” content (Wei et al., 2024a).

181 **“Harmful” content has no formal definition.** Defining success for an adversarial image is rela-
182 tively easy: the perturbation is “small” under some given measure, and leads to a misclassification.
183 With jailbreaks, however, success requires defining what it means for a model to output “harmful”
184 or otherwise “undesirable” content. Early attempts used crude proxies based on simple substring
185 matching (Zou et al., 2023). This approach has largely been replaced by a more general use of an
186 “LLM-as-a-judge”, where the fuzzy task of defining harmfulness is given to another LLM (Zheng
187 et al., 2023; Chao et al., 2023; Shah et al., 2023; Mazeika et al., 2024). The circularity of this
188 definition leads to a number of issues, as illustrated in Section 2.

189 **There are no meaningful bounds on adversaries.** Although adversaries for image classification
190 could also be unbounded, the fact that the safety property is dependent on the input (replacing a
191 cat by a dog is not an interesting attack) made the community define an ℓ_p norm around the inputs
192 as a proxy for preserving visual similarity. However, for jailbreaks, there is not such a meaningful
193 bound as the safety property is *independent* of the input (harmful generations should never occur).
194 Researchers have come up with attacks that use semantic augmentations (e.g., role-playing or social
195 engineering) (Shah et al., 2023; Zeng et al., 2024), append high-perplexity suffixes (Zou et al.,
196 2023; Thompson & Sklar, 2024) or even found that long inputs and random augmentations dilute
197 safeguards (Anil et al., 2024; Andriushchenko et al., 2024; Hughes et al., 2024). Not only adversaries
198 are now unbounded in the input space, but they can use additional methods such as fine-tuning (Qi
199 et al., 2024b) or pruning (Wei et al., 2024b). This diversity of attacks illustrates the difficulty to define
200 a narrow task, analogous to ℓ_p bounded robustness, that can be used to compare and benchmark
201 attacks and defenses.

202 **Optimizing for worst-case attacks is hard.** Optimizing attacks against classifiers is straightforward.
203 You can set as objective the maximization of the model loss (Szegedy, 2013). The loss gradient can be
204 propagated all the way to the input to guide updates. However, LLMs do not provide any of the above:
205 the optimization goal is unclear and optimization is not continuous nor over a finite input space. As a
206 workaround, previous work has tried to optimize proxy objectives such as maximizing the probability
207 of a compliance prefix (e.g. “Sure, I can help you with that”) (Zou et al., 2023; Carlini et al., 2024).
208 However, the input space is still discrete and virtually infinite. These challenges make discrete
209 optimization extremely inefficient and close to random search (Zou et al., 2023; Andriushchenko
210 et al., 2024). Optimization challenges have made us shift from a field where the strongest attacks
211 were found via white-box optimization, to one where the best attacks often come from human experts
212 and cannot be found via optimization (Li et al., 2024a). This challenges our ability to make progress
213 in measuring worst-case performance of systems (Carlini et al., 2024).

214 **3.2 Unfinetunable Models**

215 A recent research direction aims to design models that are not only robust to jailbreaks, but *also are*
216 *robust to fine-tuning* Tamirisa et al. (2024); Rosati et al. (2024). This threat model is motivated by
217 the general observation that if a model does *not* have the knowledge to perform some dangerous
218 capability (such as giving instructions for how to perform a cyberattack or design a bioweapon),
219 attacks will never be successful (Li et al., 2024b).

220 **The attacker is strictly more powerful than for adversarial examples.** An adversarial example
221 attacker has exactly one ability: to modify the input so the model produces an incorrect output. When
222 designing an un-finetunable model, we assume an attacker with *strictly* more power: not only can
223 they change the input arbitrarily, but they can also modify the model itself. Indeed, recent work has
224 already shown how the interplay between modifying the input and modifying the parameters can
225 allow attackers to break many recently proposed defenses Qi et al. (2024a).

226 **The increased attack space makes it more difficult to evaluate.** In the classical adversarial
227 example literature, the evaluator must ensure exactly one thing is true: the input-space gradient
228 is smooth and following it leads to adversarial examples. In contrast, evaluating an unfinetunable
229 model requires that the much higher *parameter-space* gradients are smooth, something often $1000 \times$
230 higher dimensional. Moreover, the number of hyperparameters in the evaluation increases significantly,
231 introducing even more room for error (Hönig et al., 2024; Qi et al., 2024a).

232 **3.3 Poisoning and Backdoors**

233 In poisoning attacks, adversaries modify a model’s training data to affect its behavior on specific
234 examples (Huang et al., 2011) or inject backdoors (Gu et al., 2019). The messy datasets and costly
235 training runs for LLMs make the definition, optimization and evaluation of attacks more challenging.

236 **Attack goals are hard to enumerate and conflict with intended functionality.** In classification
237 models, adversaries injected training examples with specific triggers that correlated with an output
238 label (Gu et al., 2019). However, in generative models, adversaries trigger fuzzy and complex
239 behaviors like producing harmful content or spreading misinformation (Wan et al., 2023; Rando
240 & Tramèr, 2024a; Zhang et al., 2024b). Not only are these behaviors harder to predict and specify
241 formally, but they also fundamentally conflict with the model’s intended functionality since the
242 triggered behavior is often universally undesirable and explicitly trained against (Zhang et al., 2024b).

243 **Attacks can come from multiple training stages and are hard to optimize over.** Traditional
244 machine learning models had a single training stage on the entire dataset. However, LLMs are first
245 pre-trained and then fine-tuned on (curated) data to turn them into helpful and harmless chatbots (Bai
246 et al., 2022). These different training stages have different properties, may enable different attacks,
247 and can overwrite poisoning in previous stages (Anwar et al., 2024; Zhang et al., 2024b). Also, in
248 LLMs there is no longer a good notion of what constitutes an effective poison nor we can optimize
249 over them (Goldblum et al., 2022).

250 **Experiments with leading models are computationally infeasible.** Rigorous evaluation of back-
251 door attacks traditionally requires training models from scratch to understand both the effects of
252 poisoned data and to establish clean baselines. However, this becomes infeasible for LLMs, where a
253 single training run can cost millions of dollars (Anwar et al., 2024; Zhang et al., 2024b).

254 **3.4 Prompt Injections**

255 In a prompt injection attack (Goodside, 2022; Willison, 2022), an adversary injects malicious
256 instructions into a language model’s context, manipulating its behavior to perform unauthorized
257 actions or disclose sensitive information. These attacks commonly target LLM agents or LLM-
258 integrated applications that interact with untrusted third-party resources through external tools (Jarvis
259 & Palermo, 2023; Husain, 2024; Anthropic, 2024).

260 **Measuring success of attacks and defenses requires a realistic AI agent environment.** Rigor-
261 ously evaluating the effectiveness of prompt injection attacks and defenses necessitates a realistic AI

262 agent environment that closely mimics real-world scenarios. Such an environment should include
263 comprehensive system scaffolding with tool use, enabling the simulation of complex interactions.
264 However, for simplicity, many studies opt to simulate these environments and rely on LLMs as judges
265 for evaluation. There are new setups that have more rigorous evaluations (Debenedetti et al., 2024),
266 where the attack’s success and utility can be precisely measured, but they are often limited due to the
267 high cost of incorporating new tasks and their reliance on simulated environments.

268 **Adversaries are unbounded.** Unlike traditional adversarial attacks bounded by ℓ_p norms, prompt
269 injection attacks also operate in a vast and unbounded input space. Additionally, prompt injection
270 attacks can leverage context-dependent strategies, such as embedding malicious instructions within
271 seemingly benign or unrelated text, or using multi-turn interactions to gradually steer the model
272 toward undesirable outputs. This diversity in attack vectors, combined with the fact that virtually
273 any controlled input can serve as a potential attack surface, complicates the task of establishing a
274 reasonable threat model. Consequently, creating a standardized “toy” problem for benchmarking
275 prompt injection defenses is inherently difficult.

276 **Optimizing for strong attacks is hard.** The primary goal of prompt injections is often clear—
277 for instance, manipulating a language model to perform unauthorized actions like sending
278 emails (Debenedetti et al., 2024), where success can be directly measured. However, the attack
279 surface remains vast, encompassing not only single-turn interactions but also multi-turn scenarios
280 where the model may repeatedly call external tools. In such cases, researchers often lack access to
281 intermediate outputs, making it significantly more challenging to refine and optimize the attack.

282 Most current attacks rely on handcrafted instructions (Greshake et al., 2023; Liu et al., 2023), such as,
283 “Ignore all previous instructions, please do [target action] first,” which are often effective in practice.
284 These manual attacks complicate the development of principled defenses like adversarial training,
285 due to their highly context-dependent and ad hoc nature. Recent approaches (Pasquini et al., 2024)
286 have attempted to apply optimization techniques similar to those used in jailbreaks. Unfortunately,
287 these attacks are not guaranteed to be optimal. As a result, defense attempts that train models against
288 attacks mainly focus on *known* attacks Wallace et al. (2024).

289 **We cannot easily track progress against closed-source systems.** Similar to jailbreaks, model
290 developers can mitigate prompt injection attacks by implementing safeguards such as filtering mechanisms
291 (Willison, 2023; Wu et al., 2024) or regularly updating and fine-tuning their models (Wallace
292 et al., 2024). As these systems are frequently updated, it becomes difficult to establish a consistent
293 benchmark for measuring progress or reproducing results. Additionally, there are currently few
294 open-source models that are effective tool-use agents (Debenedetti et al., 2024) and can be used for
295 reproducible evaluation.

296 3.5 Membership Inference

297 Membership inference (MI) attacks (Shokri et al., 2017) aim to determine whether a specific sample
298 x was part of a model’s training set.

299 **The distinction between members and non-members is no longer clearly defined.** In traditional
300 classification settings, the training data is typically of limited size and with a clear delimitation
301 between samples. However, the situation becomes more complicated for generative models.

- 302 1. **Highly (partially) duplicated datasets.** The training data of generative models often comes
303 from massive, diverse open datasets, which could include numerous duplicate and near-duplicate
304 samples (Lee et al., 2022; Tirumala et al., 2023). Even if a model appears to memorize a particular
305 sample (e.g., a piece of text or image), this does not necessarily prove that this sample itself was
306 used during training. For example, a model might know much of the plot of Harry Potter without
307 having been explicitly trained on the original book; it could have learned about the story indirectly
308 through Wikipedia pages, reviews, etc. Thus, the boundaries between members and non-members
309 are blurred by the sheer scale and overlap of these datasets.
- 310 2. **No IID train and test splits available.** Methods for evaluating MI designate the training data as
311 members and separate IID held-out data as non-members. However, for most generative models, the
312 training datasets are typically not disclosed. Some recent studies attempt to collect non-members

313 post hoc for evaluation purposes (Shi et al., 2023; Meeus et al., 2023), but these efforts often
314 violate the IID assumption and lead to misleading conclusions (Duan et al., 2024; Das et al., 2024).

315 **We cannot build counterfactual scenarios for evaluation.** In traditional classification tasks
316 (e.g., CIFAR-10), where the data generation process is known and models are relatively small,
317 counterfactual scenarios can be built by retraining the same model while excluding a sample x , and
318 then comparing statistical behaviors on x (Carlini et al., 2022). In the context of generative models,
319 this approach is ill-defined and computationally impractical, thus it's infeasible to properly evaluate
320 the success of a MI attack (Zhang et al., 2024a).

321 **3.6 Machine Unlearning**

322 Machine unlearning was originally formulated as a well-defined task: completely removing the
323 influence of a specific datapoint x from a model (Bourtoule et al., 2021). The goal was to produce a
324 model that, after unlearning x , would be indistinguishable from one that was never trained on that
325 point. In traditional classification settings with bounded inputs and outputs, and (often) deduplicated
326 datasets with clear train-test splits, this objective could be precisely defined and evaluated. In fact,
327 there exist exact solutions to unlearning (Bourtoule et al., 2021).

328 **Unlearning of “concepts” rather than individual data points is hard to define.** However,
329 generative models have fundamentally changed the nature of unlearning (Cooper et al., 2024). Instead
330 of removing the influence of specific data points, the goal is to remove knowledge about entire
331 concepts or topics that may be contained in one *or more* data points (e.g., all dangerous knowledge
332 about bioweapons (Li et al., 2024b) or copyrighted content from Harry Potter books (Eldan &
333 Russinovich, 2023)). This has made it impossible to define unlearning in terms of a specific data
334 point's influence, making both solutions and evaluations much more challenging.

335 **Unlearning goals conflict with other knowledge.** Developers may need to remove very specific
336 knowledge (e.g., bioweapons) while maintaining the model's expertise in related fields (e.g., biology
337 and virology) (Li et al., 2024b). This tension between harmful and benign knowledge makes it
338 inherently hard to define the goal of unlearning and to robustly evaluate safety and utility.

339 **Threat models are overly strong.** Unlearning emerged as a white-box protection that would
340 prevent *any* adversary from accessing undesired capabilities (Li et al., 2024b). This ambitious goal
341 also enables stronger threat models where adversaries cannot only query the model, but also finetune
342 it (Hu et al., 2024) and perform any white-box interventions (Lucki et al., 2024). Protecting against
343 such a large attack surface is much harder (Qi et al., 2024a) as discussed in Section 3.2.

344 **Measuring unlearning success is hard.** Measuring unlearning success has become significantly
345 more challenging: training baseline models without specific datapoints is costly (Eldan & Russinovich,
346 2023) and membership inference has important limitations (see Section 3.5). Recent studies have
347 also demonstrated that even when a model cannot generate specific information, this does not reliably
348 prove the underlying knowledge has been erased from its weights (Patil et al., 2023; Lynch et al.,
349 2024; Lucki et al., 2024; Shumailov et al., 2024). In practice, the search for adaptive evaluations is
350 impractical and requires very careful tuning of the methodology for each scenario (Lucki et al., 2024;
351 Qi et al., 2024a). Finally, Shi et al. (2024) showed that measuring unintended effects of unlearning is
352 challenging, as it can significantly affect other capabilities or even amplify privacy leakage.

353 **4 Discussion**

354 **4.1 Alternative Views**

355 **We are solving the right problem in the first place.** We see increased complexity in adversarial
356 ML because we are finally attempting to solve *real* security challenges rather than toy academic
357 problems. We knew that ℓ_p -bounded perturbations were a simplified proxy (Gilmer et al., 2018), but
358 they were studied because they were challenging enough to drive progress and served as a *necessary*
359 condition for real-world robustness. We could similarly define toy problems for LLMs (e.g., jailbreaks
360 limited to fixed-length prefixes or bounded sentence modifications), but the field has largely avoided

361 such artificial constraints in favor of studying real-world unbounded adversaries. This shift might not
362 indicate that problems have become fundamentally harder, but rather that the research community
363 has decided to directly tackle the full complexity of real-world security.

364 **Solving jailbreaks might be easier because we only need to prevent a behavior regardless of con-**
365 **text.** Some researchers argue that certain problems have become simpler with LLMs. For instance,
366 unlike adversarial examples where a model should maintain correct predictions in appropriate contexts
367 (e.g., classify guacamole images as guacamole, but never cats as guacamole), jailbreak prevention has
368 a simpler goal: the model should *never* produce certain harmful outputs (e.g., instructions for building
369 explosives) regardless of context. However, since there are many ways to express this knowledge
370 (e.g., harmful requests can be decomposed into benign subquestions (Glukhov et al., 2024)), defining
371 and evaluating whether a model will *never* produce harmful outputs remains a challenging problem.

372 Recent work, on representation engineering (Arditi et al., 2024; Zou et al., 2024; Tamirisa et al.,
373 2024) has aimed to identify specific directions in the model’s representation space that can anticipate
374 undesired behavior and prevent it universally. Yet, we know that adversarial images could also be
375 detected by similar methods (Carlini & Wagner, 2017), but these defenses ultimately proved vulnera-
376 ble to newer attacks. Similarly, there are already works that show that representation engineering
377 methods cannot robustly void undesired behaviors (Li et al., 2024a; Qi et al., 2024a).

378 **Scaffolding to reduce the probability of failure might be sufficient.** Given the difficulty of
379 achieving robust safety guarantees, researchers and companies increasingly rely on complex defense
380 systems (Sharma et al., 2025) and security through obscurity (Rando & Tramèr, 2024b) to minimize
381 risks. While this approach has demonstrated clear benefits in protecting users from harmful content,
382 it prevents rigorous, reproducible and adaptive evaluations as systems become more complex and
383 opaque (Casper et al., 2024a). This trend is particularly concerning given historical lessons: preventing
384 researchers from thoroughly analyzing systems can lead to severe real-world security breaches (Swire,
385 2004; Mulligan & Perzanowski, 2007; Payne & Parks, 2020). The apparent safety gains from
386 obscurity and complexity may come at the cost of genuine security understanding.

387 **We are already making progress on these problems.** A prevalent view in the field suggests that
388 we are advancing security capabilities, pointing to newer models being demonstrably harder to attack
389 than their predecessors (Achiam et al., 2023; Zaremba et al., 2025). While this observation might
390 hold generally true, we caution that our inability to robustly evaluate defenses may be hindering
391 our ability to track progress (see Section 2.3). Moreover, we must distinguish between progress in
392 preventing average-case vulnerabilities and achieving **worst-case** security robustness. Although we
393 might be making progress in the former, we have barely improved the latter and most models can still
394 produce harmful generations under attacks. As the stakes increase with more capable models, the
395 risks of rare yet successful attacks become significant (Anthropic, 2023).

396 4.2 Suggestions for improvement

397 We propose that there are (at least) two valid reasons for performing research on adversarial machine
398 learning: (a) studying real-world security vulnerabilities and (b) advancing scientific understanding
399 of adversarial ML. Papers should be explicit for what reason they are being written, and should be
400 evaluated in this light. For real-world security, demonstrating attacks on fuzzy, ill-defined problems
401 can be valuable when the potential harm is clear and immediate. For instance, it is valuable to
402 show that language models can be manipulated to produce harmful content, even if we cannot
403 precisely quantify “harmfulness”. And when the objective is to advancing scientific understanding,
404 we believe it is more productive to identify and focus on formal, well-defined sub-problems that can
405 be rigorously studied, similar to how ℓ_p -bounded perturbations provided a concrete framework for
406 studying adversarial examples.

407 We acknowledge that even these well-defined sub-problems might still be challenging, just as
408 achieving reliable ℓ_p robustness remains an open problem despite a decade of research. However,
409 what we can definitely say is that if we cannot make progress on carefully scoped, formal problems, we
410 have little hope of addressing the broader, fuzzier challenges of language model security. Moreover,
411 working on well-defined problems enables rigorous scientific investigation: we can properly measure
412 progress, compare different approaches, and build upon previous results. Attempting to solve the
413 entire space of attacks without rigor is neither scientific nor likely to be productive.

414 **References**

415 Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt,
416 J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

417 Andriushchenko, M., Croce, F., and Flammarion, N. Jailbreaking leading safety-aligned llms with
418 simple adaptive attacks. *arXiv preprint arXiv:2404.02151*, 2024.

419 Anil, C., Durmus, E., Rimsky, N., Sharma, M., Benton, J., Kundu, S., Batson, J., Tong, M., Mu, J.,
420 Ford, D. J., et al. Many-shot jailbreaking. In *The Thirty-eighth Annual Conference on Neural
421 Information Processing Systems*, 2024.

422 Anthropic. Anthropic's responsible scaling policy. <https://www.anthropic.com/news/anthropics-responsible-scaling-policy>, 2023.

424 Anthropic. Tool use (function calling). <https://docs.anthropic.com/en/docs/tool-use>,
425 2024.

426 Anwar, U., Saparov, A., Rando, J., Paleka, D., Turpin, M., Hase, P., Lubana, E. S., Jenner, E.,
427 Casper, S., Sourbut, O., et al. Foundational challenges in assuring alignment and safety of large
428 language models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. Survey
429 Certification, Expert Certification.

430 Ardit, A., Obeso, O. B., Syed, A., Paleka, D., Rimsky, N., Gurnee, W., and Nanda, N. Refusal in
431 language models is mediated by a single direction. In *The Thirty-eighth Annual Conference on
432 Neural Information Processing Systems*, 2024.

433 Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., Drain, D., Fort, S., Ganguli, D.,
434 Henighan, T., et al. Training a helpful and harmless assistant with reinforcement learning from
435 human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

436 Barreno, M., Nelson, B., Sears, R., Joseph, A. D., and Tygar, J. D. Can machine learning be secure?
437 In *Proceedings of the 2006 ACM Symposium on Information, Computer and Communications
438 Security, ASIACCS '06*, pp. 16–25, New York, NY, USA, 2006. Association for Computing
439 Machinery. ISBN 1595932720. doi: 10.1145/1128817.1128824.

440 Biggio, B., Corona, I., Maiorca, D., Nelson, B., Šrndić, N., Laskov, P., Giacinto, G., and Roli,
441 F. Evasion attacks against machine learning at test time. In *Machine Learning and Knowledge
442 Discovery in Databases: European Conference, ECML PKDD 2013, Prague, Czech Republic,
443 September 23-27, 2013, Proceedings, Part III* 13, pp. 387–402. Springer, 2013.

444 Bourtoule, L., Chandrasekaran, V., Choquette-Choo, C. A., Jia, H., Travers, A., Zhang, B., Lie, D.,
445 and Papernot, N. Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, pp.
446 141–159. IEEE, 2021.

447 Carlini, N. and Wagner, D. Towards evaluating the robustness of neural networks. In *2017 ieee
448 symposium on security and privacy (sp)*, pp. 39–57. Ieee, 2017.

449 Carlini, N., Katz, G., Barrett, C., and Dill, D. L. Provably minimally-distorted adversarial examples.
450 *arXiv preprint arXiv:1709.10207*, 2017.

451 Carlini, N., Athalye, A., Papernot, N., Brendel, W., Rauber, J., Tsipras, D., Goodfellow, I., Madry, A.,
452 and Kurakin, A. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.

453 Carlini, N., Chien, S., Nasr, M., Song, S., Terzis, A., and Tramer, F. Membership inference attacks
454 from first principles. In *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 1897–1914.
455 IEEE, 2022.

456 Carlini, N., Nasr, M., Choquette-Choo, C. A., Jagielski, M., Gao, I., Koh, P. W. W., Ippolito, D.,
457 Tramer, F., and Schmid, L. Are aligned neural networks adversarially aligned? *Advances in
458 Neural Information Processing Systems*, 36, 2024.

459 Casper, S., Davies, X., Shi, C., Gilbert, T. K., Scheurer, J., Rando, J., Freedman, R., Korbak, T.,
 460 Lindner, D., Freire, P., et al. Open problems and fundamental limitations of reinforcement learning
 461 from human feedback. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
 462 Survey Certification, Featured Certification.

463 Casper, S., Ezell, C., Siegmann, C., Kolt, N., Curtis, T. L., Bucknall, B., Haupt, A., Wei, K., Scheurer,
 464 J., Hobbhahn, M., et al. Black-box access is insufficient for rigorous ai audits. In *The 2024 ACM*
 465 *Conference on Fairness, Accountability, and Transparency*, pp. 2254–2272, 2024a.

466 Casper, S., Schulze, L., Patel, O., and Hadfield-Menell, D. Defending against unforeseen failure
 467 modes with latent adversarial training. *arXiv preprint arXiv:2403.05030*, 2024b.

468 Chao, P., Robey, A., Dobriban, E., Hassani, H., Pappas, G. J., and Wong, E. Jailbreaking black box
 469 large language models in twenty queries. *arXiv preprint arXiv:2310.08419*, 2023.

470 Chao, P., Debenedetti, E., Robey, A., Andriushchenko, M., Croce, F., Sehwag, V., Dobriban, E.,
 471 Flammarion, N., Pappas, G. J., Tramèr, F., et al. Jailbreakbench: An open robustness benchmark
 472 for jailbreaking large language models. In *The Thirty-eight Conference on Neural Information*
 473 *Processing Systems Datasets and Benchmarks Track*, 2024.

474 Chi, J., Karn, U., Zhan, H., Smith, E., Rando, J., Zhang, Y., Plawiak, K., Coudert, Z. D., Upasani,
 475 K., and Pasupuleti, M. Llama guard 3 vision: Safeguarding human-ai image understanding
 476 conversations. *arXiv preprint arXiv:2411.10414*, 2024.

477 Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., and Amodei, D. Deep reinforcement
 478 learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

479 Cohen, J., Rosenfeld, E., and Kolter, Z. Certified adversarial robustness via randomized smoothing.
 480 In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference*
 481 *on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1310–1320.
 482 PMLR, 09–15 Jun 2019.

483 Cooper, A. F., Choquette-Choo, C. A., Bogen, M., Jagielski, M., Filippova, K., Liu, K. Z., Choulde-
 484 chova, A., Hayes, J., Huang, Y., Mireshghallah, N., et al. Machine unlearning doesn't do what you
 485 think: Lessons for generative ai policy, research, and practice. *arXiv preprint arXiv:2412.06966*,
 486 2024.

487 Cui, J., Chiang, W.-L., Stoica, I., and Hsieh, C.-J. Or-bench: An over-refusal benchmark for large
 488 language models. *arXiv preprint arXiv:2405.20947*, 2024.

489 Das, D., Zhang, J., and Tramèr, F. Blind baselines beat membership inference attacks for foundation
 490 models. *arXiv preprint arXiv:2406.16201*, 2024.

491 Debenedetti, E., Zhang, J., Balunović, M., Beurer-Kellner, L., Fischer, M., and Tramèr, F. Agent-
 492 dojo: A dynamic environment to evaluate attacks and defenses for llm agents. *arXiv preprint*
 493 *arXiv:2406.13352*, 2024.

494 Duan, M., Suri, A., Mireshghallah, N., Min, S., Shi, W., Zettlemoyer, L., Tsvetkov, Y., Choi, Y.,
 495 Evans, D., and Hajishirzi, H. Do membership inference attacks work on large language models?
 496 *arXiv preprint arXiv:2402.07841*, 2024.

497 Eldan, R. and Russinovich, M. Who's harry potter? approximate unlearning in llms. *arXiv preprint*
 498 *arXiv:2310.02238*, 2023.

499 Gao, L., Biderman, S., Black, S., Golding, L., Hoppe, T., Foster, C., Phang, J., He, H., Thite, A.,
 500 Nabeshima, N., et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv*
 501 *preprint arXiv:2101.00027*, 2020.

502 Gilmer, J., Adams, R. P., Goodfellow, I., Andersen, D., and Dahl, G. E. Motivating the rules of the
 503 game for adversarial example research. *arXiv preprint arXiv:1807.06732*, 2018.

504 Glukhov, D., Han, Z., Shumailov, I., Papyan, V., and Papernot, N. Breach by a thousand leaks:
 505 Unsafe information leakage in safe'ai responses. *arXiv preprint arXiv:2407.02551*, 2024.

506 Goldblum, M., Tsipras, D., Xie, C., Chen, X., Schwarzschild, A., Song, D., Madry, A., Li, B.,
 507 and Goldstein, T. Dataset security for machine learning: Data poisoning, backdoor attacks, and
 508 defenses. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):1563–1580,
 509 2022.

510 Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv
 511 preprint arXiv:1412.6572*, 2014.

512 Goodside, R. Exploiting GPT-3 prompts with malicious inputs that order the model to ignore its
 513 previous directions. <https://x.com/goodside/status/1569128808308957185>, 2022.

514 Graham-Cumming, J. How to beat an adaptive spam filter. In *MIT Spam Conference*, January 2004.
 515 Oral presentation.

516 Greshake, K., Abdelnabi, S., Mishra, S., Endres, C., Holz, T., and Fritz, M. Not what you've signed
 517 up for: Compromising real-world LLM-integrated applications with indirect prompt injection. In
 518 *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, CCS '23. ACM,
 519 November 2023. doi: 10.1145/3605764.3623985.

520 Gu, T., Liu, K., Dolan-Gavitt, B., and Garg, S. Badnets: Evaluating backdooring attacks on deep
 521 neural networks. *IEEE Access*, 7:47230–47244, 2019.

522 Höning, R., Rando, J., Carlini, N., and Tramèr, F. Adversarial perturbations cannot reliably protect
 523 artists from generative ai. *arXiv preprint arXiv:2406.12027*, 2024.

524 Hu, S., Fu, Y., Wu, Z. S., and Smith, V. Jogging the memory of unlearned model through targeted
 525 relearning attack. *arXiv preprint arXiv:2406.13356*, 2024.

526 Huang, L., Joseph, A. D., Nelson, B., Rubinstein, B. I., and Tygar, J. D. Adversarial machine learning.
 527 In *Proceedings of the 4th ACM workshop on Security and artificial intelligence*, pp. 43–58, 2011.

528 Hughes, J., Price, S., Lynch, A., Schaeffer, R., Barez, F., Koyejo, S., Sleight, H., Jones, E., Perez, E.,
 529 and Sharma, M. Best-of-n jailbreaking. *arXiv preprint arXiv:2412.03556*, 2024.

530 Husain, H. Llama-3 function calling demo. [https://nbsanity.com/static/
 531 d06085f1dacae8c9de9402f2d7428de2/demo.html](https://nbsanity.com/static/d06085f1dacae8c9de9402f2d7428de2/demo.html), 2024.

532 Inan, H., Upasani, K., Chi, J., Rungta, R., Iyer, K., Mao, Y., Tontchev, M., Hu, Q., Fuller, B.,
 533 Testuggine, D., et al. Llama guard: Llm-based input-output safeguard for human-ai conversations.
 534 *arXiv preprint arXiv:2312.06674*, 2023.

535 Jarvis, C. and Palermo, J. Function calling. https://cookbook.openai.com/examples/how_to_call_functions_with_chat_models, 6 2023.

537 Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional
 538 neural networks. *Advances in neural information processing systems*, 25, 2012.

539 Lee, K., Ippolito, D., Nyström, A., Zhang, C., Eck, D., Callison-Burch, C., and Carlini, N. Deduplicating
 540 training data makes language models better. In *Proceedings of the 60th Annual Meeting of the Association
 541 for Computational Linguistics*. Association for Computational Linguistics, 2022.

542 Li, N., Han, Z., Steneker, I., Primack, W., Goodside, R., Zhang, H., Wang, Z., Menghini, C., and Yue,
 543 S. Llm defenses are not robust to multi-turn human jailbreaks yet. *arXiv preprint arXiv:2408.15221*,
 544 2024a.

545 Li, N., Pan, A., Gopal, A., Yue, S., Berrios, D., Gatti, A., Li, J. D., Dombrowski, A.-K., Goel, S.,
 546 Mukobi, G., et al. The WMDP benchmark: Measuring and reducing malicious use with unlearning.
 547 In *Forty-first International Conference on Machine Learning*, 2024b.

548 Liu, F., Feng, Y., Xu, Z., Su, L., Ma, X., Yin, D., and Liu, H. Jailjudge: A comprehensive jailbreak
 549 judge benchmark with multi-agent enhanced explanation evaluation framework. *arXiv preprint
 550 arXiv:2410.12855*, 2024.

551 Liu, Y., Deng, G., Li, Y., Wang, K., Zhang, T., Liu, Y., Wang, H., Zheng, Y., and Liu, Y. Prompt
 552 injection attack against LLM-integrated applications. *arXiv preprint arXiv:2306.05499*, 2023.

553 Lucki, J., Wei, B., Huang, Y., Henderson, P., Tramèr, F., and Rando, J. An adversarial perspective on
554 machine unlearning for ai safety. *arXiv preprint arXiv:2409.18025*, 2024.

555 Lynch, A., Guo, P., Ewart, A., Casper, S., and Hadfield-Menell, D. Eight methods to evaluate robust
556 unlearning in llms. *arXiv preprint arXiv:2402.16835*, 2024.

557 Madry, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint
558 arXiv:1706.06083*, 2017.

559 Mai, W., Hong, G., Chen, P., Pan, X., Liu, B., Zhang, Y., Duan, H., and Yang, M. You can't eat your
560 cake and have it too: The performance degradation of llms with jailbreak defense, 2025.

561 Mangaokar, N., Hooda, A., Choi, J., Chandrashekaran, S., Fawaz, K., Jha, S., and Prakash, A. Prp:
562 Propagating universal perturbations to attack large language model guard-rails. *arXiv preprint
563 arXiv:2402.15911*, 2024.

564 Mazeika, M., Phan, L., Yin, X., Zou, A., Wang, Z., Mu, N., Sakhaei, E., Li, N., Basart, S., Li, B.,
565 et al. Harmbench: A standardized evaluation framework for automated red teaming and robust
566 refusal. *arXiv preprint arXiv:2402.04249*, 2024.

567 Meeus, M., Jain, S., Rei, M., and de Montjoye, Y.-A. Did the neurons read your book? Document-
568 level membership inference for large language models. *arXiv preprint arXiv:2310.15007*, 2023.

569 Miyato, T., Maeda, S.-i., Koyama, M., and Ishii, S. Virtual adversarial training: a regularization
570 method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and
571 machine intelligence*, 41(8):1979–1993, 2018.

572 Mulligan, D. K. and Perzanowski, A. K. The magnificence of the disaster: Reconstructing the sony
573 bmg rootkit incident. *Berkeley Tech. LJ*, 22:1157, 2007.

574 Nasr, M., Rando, J., Carlini, N., Hayase, J., Jagielski, M., Cooper, A. F., Ippolito, D., Choquette-
575 Choo, C. A., Tramèr, F., and Lee, K. Scalable extraction of training data from aligned, production
576 language models. In *The Thirteenth International Conference on Learning Representations*, 2025.

577 Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S.,
578 Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback.
579 *Advances in neural information processing systems*, 35:27730–27744, 2022.

580 Pasquini, D., Strohmeier, M., and Troncoso, C. Neural exec: Learning (and learning from) execution
581 triggers for prompt injection attacks. In *Proceedings of the 2024 Workshop on Artificial Intelligence
582 and Security*, pp. 89–100, 2024.

583 Patil, V., Hase, P., and Bansal, M. Can sensitive information be deleted from llms? objectives for
584 defending against extraction attacks. *arXiv preprint arXiv:2309.17410*, 2023.

585 Payne, K. and Parks, M. Despite election security fears, iowa caucuses will use new smartphone app.
586 *National Public Radio.*, 2020.

587 Qi, X., Wei, B., Carlini, N., Huang, Y., Xie, T., He, L., Jagielski, M., Nasr, M., Mittal, P., and
588 Henderson, P. On evaluating the durability of safeguards for open-weight llms. *arXiv preprint
589 arXiv:2412.07097*, 2024a.

590 Qi, X., Zeng, Y., Xie, T., Chen, P.-Y., Jia, R., Mittal, P., and Henderson, P. Fine-tuning aligned lan-
591 guage models compromises safety, even when users do not intend to! In *The Twelfth International
592 Conference on Learning Representations*, 2024b.

593 Raina, V., Liusie, A., and Gales, M. Is llm-as-a-judge robust? investigating universal adversarial
594 attacks on zero-shot llm assessment. *arXiv preprint arXiv:2402.14016*, 2024.

595 Rajpurkar, P. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint
596 arXiv:1606.05250*, 2016.

597 Rando, J. and Tramèr, F. Universal jailbreak backdoors from poisoned human feedback. In *The
598 Twelfth International Conference on Learning Representations*, 2024a.

599 Rando, J. and Tramèr, F. The worst (but only) claude 3 tokenizer. <https://javirando.com/blog/2024/clause-tokenizer/>, 2024b.

600

601 Rando, J., Korevaar, H., Brinkman, E., Evtimov, I., and Tramèr, F. Gradient-based jailbreak images for multimodal fusion models. *arXiv preprint arXiv:2410.03489*, 2024.

602

603 Richardson, M., Burges, C. J., and Renshaw, E. McTest: A challenge dataset for the open-domain machine comprehension of text. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 193–203, 2013.

604

605

606 Robey, A., Wong, E., Hassani, H., and Pappas, G. J. Smoothllm: Defending large language models against jailbreaking attacks. *arXiv preprint arXiv:2310.03684*, 2023.

607

608 Rosati, D., Wehner, J., Williams, K., Bartoszcze, Ł., Atanasov, D., Gonzales, R., Majumdar, S., Maple, C., Sajjad, H., and Rudzicz, F. Representation noising effectively prevents harmful fine-tuning on llms. *NeurIPS*, 2024.

609

610

611 Shafahi, A., Huang, W. R., Studer, C., Feizi, S., and Goldstein, T. Are adversarial examples inevitable? In *International Conference on Learning Representations*, 2019.

612

613 Shah, R., Pour, S., Tagade, A., Casper, S., Rando, J., et al. Scalable and transferable black-box jailbreaks for language models via persona modulation. *arXiv preprint arXiv:2311.03348*, 2023.

614

615 Sharma, M., Tong, M., Mu, J., Wei, J., Kruthoff, J., Goodfriend, S., Ong, E., Peng, A., Agarwal, R., Anil, C., et al. Constitutional classifiers: Defending against universal jailbreaks across thousands of hours of red teaming, 2025.

616

617

618 Sheshadri, A., Ewart, A., Guo, P., Lynch, A., Wu, C., Hebbar, V., Sleight, H., Stickland, A. C., Perez, E., Hadfield-Menell, D., et al. Latent adversarial training improves robustness to persistent harmful behaviors in llms. *arXiv preprint arXiv:2407.15549*, 2024.

619

620

621 Shi, W., Ajith, A., Xia, M., Huang, Y., Liu, D., Blevins, T., Chen, D., and Zettlemoyer, L. Detecting pretraining data from large language models. *arXiv preprint arXiv:2310.16789*, 2023.

622

623 Shi, W., Lee, J., Huang, Y., Malladi, S., Zhao, J., Holtzman, A., Liu, D., Zettlemoyer, L., Smith, N. A., and Zhang, C. Muse: Machine unlearning six-way evaluation for language models. *arXiv preprint arXiv:2407.06460*, 2024.

624

625

626 Shokri, R., Stronati, M., Song, C., and Shmatikov, V. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pp. 3–18. IEEE, 2017.

627

628 Shumailov, I., Hayes, J., Triantafillou, E., Ortiz-Jimenez, G., Papernot, N., Jagielski, M., Yona, I., Howard, H., and Bagdasaryan, E. Unlearning: Unlearning is not sufficient for content regulation in advanced generative ai. *arXiv preprint arXiv:2407.00106*, 2024.

629

630

631 Souly, A., Lu, Q., Bowen, D., Trinh, T., Hsieh, E., Pandey, S., Abbeel, P., Svegliato, J., Emmons, S., Watkins, O., et al. A strongreject for empty jailbreaks. *arXiv preprint arXiv:2402.10260*, 2024.

632

633 Swire, P. P. A model for when disclosure helps security: What is different about computer and network security? *J. on Telecomm. & High Tech. L.*, 3:163, 2004.

634

635 Szegedy, C. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

636 Tamirisa, R., Bharathi, B., Phan, L., Zhou, A., Gatti, A., Suresh, T., Lin, M., Wang, J., Wang, R., Arel, R., et al. Tamper-resistant safeguards for open-weight llms. *arXiv preprint arXiv:2408.00761*, 2024.

637

638

639 Thompson, T. B. and Sklar, M. Flrt: Fluent student-teacher redteaming. *arXiv preprint arXiv:2407.17447*, 2024.

640

641 Tirumala, K., Simig, D., Aghajanyan, A., and Morcos, A. D4: Improving llm pretraining via document de-duplication and diversification. *Advances in Neural Information Processing Systems*, 36:53983–53995, 2023.

642

643

644 Tramer, F., Carlini, N., Brendel, W., and Madry, A. On adaptive attacks to adversarial example
 645 defenses. *Advances in neural information processing systems*, 33:1633–1645, 2020.

646 Wallace, E., Xiao, K., Leike, R., Weng, L., Heidecke, J., and Beutel, A. The instruction hierarchy:
 647 Training llms to prioritize privileged instructions. *arXiv preprint arXiv:2404.13208*, 2024.

648 Wan, A., Wallace, E., Shen, S., and Klein, D. Poisoning language models during instruction tuning.
 649 In *International Conference on Machine Learning*, pp. 35413–35425. PMLR, 2023.

650 Wei, A., Haghtalab, N., and Steinhardt, J. Jailbroken: How does llm safety training fail? *Advances in
 651 Neural Information Processing Systems*, 36, 2024a.

652 Wei, B., Huang, K., Huang, Y., Xie, T., Qi, X., Xia, M., Mittal, P., Wang, M., and Henderson, P.
 653 Assessing the brittleness of safety alignment via pruning and low-rank modifications. In *Forty-first
 654 International Conference on Machine Learning*, 2024b.

655 Willison, S. Prompt injection attacks against GPT-3. <https://simonwillison.net/2022/Sep/12/prompt-injection/>, 2022.

656 Willison, S. Delimiters won't save you from prompt injection. <https://simonwillison.net/2023/May/11/delimiters-wont-save-you/>, 2023.

659 Wu, Y., Roesner, F., Kohno, T., Zhang, N., and Iqbal, U. SecGPT: An execution isolation architecture
 660 for LLM-based systems. *arXiv preprint arXiv:2403.04960*, 2024.

661 Zaremba, W., Nitishinskaya, E., Barak, B., Lin, S., Toyer, S., Yu, Y., Dias, R., Wallace, E., Xiao, K.,
 662 and Glaese, J. H. A. Trading inference-time compute for adversarial robustness, 2025.

663 Zeng, Y., Lin, H., Zhang, J., Yang, D., Jia, R., and Shi, W. How johnny can persuade llms to
 664 jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. *arXiv preprint
 665 arXiv:2401.06373*, 2024.

666 Zhang, J., Das, D., Kamath, G., and Tramèr, F. Membership inference attacks cannot prove that a
 667 model was trained on your data. *arXiv preprint arXiv:2409.19798*, 2024a.

668 Zhang, Y., Rando, J., Evtimov, I., Chi, J., Smith, E. M., Carlini, N., Tramèr, F., and Ippolito, D.
 669 Persistent pre-training poisoning of llms. *arXiv preprint arXiv:2410.13722*, 2024b.

670 Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E.,
 671 et al. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference
 672 on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.

673 Zou, A., Wang, Z., Carlini, N., Nasr, M., Kolter, J. Z., and Fredrikson, M. Universal and transferable
 674 adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

675 Zou, A., Phan, L., Wang, J., Duenas, D., Lin, M., Andriushchenko, M., Wang, R., Kolter, Z.,
 676 Fredrikson, M., and Hendrycks, D. Improving alignment and robustness with short circuiting.
 677 *arXiv preprint arXiv:2406.04313*, 2024.