

Red-Teaming for Generative AI: Silver Bullet or Security Theater?

Michael Feffer, Anusha Sinha, Zachary C. Lipton, Hoda Heidari

Carnegie Mellon University
mfeffer@andrew.cmu.edu, asinha@sei.cmu.edu,
{zlipton, hheidari}@andrew.cmu.edu

Abstract

In response to rising concerns surrounding the safety, security, and trustworthiness of Generative AI (GenAI) models, practitioners and regulators alike have pointed to *AI red-teaming* as a key component of their strategies for identifying and mitigating these risks. However, despite AI red-teaming’s central role in policy discussions and corporate messaging, significant questions remain about what precisely it means, what role it can play in regulation, and how it relates to conventional red-teaming practices as originally conceived in the field of cybersecurity. In this work, we identify recent cases of red-teaming activities in the AI industry and conduct an extensive survey of the relevant research literature to characterize the scope, structure, and criteria for AI red-teaming practices. Our analysis reveals that prior methods and practices of AI red-teaming diverge along several axes, including the purpose of the activity (which is often vague), the artifact under evaluation, the setting in which the activity is conducted (e.g., actors, resources, and methods), and the resulting decisions it informs (e.g., reporting, disclosure, and mitigation). In light of our findings, we argue that while red-teaming may be a valuable big-tent idea for characterizing a broad set of activities and attitudes aimed at improving the behavior of GenAI models, gestures towards red-teaming as a panacea for every possible risk verge on *security theater*. To move toward a more robust toolbox of evaluations for generative AI, we synthesize our recommendations into a question bank meant to guide and scaffold future AI red-teaming practices.

1 Introduction

In recent years, generative AI technologies, including large language models (LLMs) [134, 94] image and video generation models [104, 110], and audio generation models [44, 8] have captured the public imagination. While many view the proliferation and accessibility of these tools favorably, envisioning a boon to productivity, creativity, and economic growth, concerns have emerged that the rapid adoption of these powerful models might unleash new categories of societal harms. These concerns have gained credibility owing to several well-publicized problematic incidents where such AI output text expressing discriminatory sentiment towards marginalized groups [86, 54, 93], created images reflecting harmful stereotypes [80], and enabled the generation of deep-fake audio in a fashion that has been likened to *digital blackface* [47]. These issues are compounded by the lack of transparency and accountability surrounding the development and evaluation of these models [94, 145, 18].

In answer to the mounting worry over the safety, security, and trustworthiness of generative AI models, practitioners and policy-makers alike have pointed to *red-teaming* as an integral part of their strategies to identify and address related risks, with the goal of ensuring some notion of alignment with human and societal values [11, 87, 21]. Notably, the US presidential Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence [59] mentions red-teaming eight times, defining the practice as follows:

“The term ‘AI red-teaming’ means a structured testing effort to find flaws and vulnerabilities in an AI system, often in a controlled environment and in collaboration with developers of AI. Artificial Intelligence red-teaming is most often performed by dedicated ‘red teams’ that adopt adversarial methods to identify flaws and vulnerabilities, such as harmful or discriminatory outputs from an AI system, unforeseen or undesirable system behaviors, limitations, or potential risks associated with the misuse of the system.”

The order mandates the Secretary of Commerce and other federal agencies to develop guidelines, standards, and best practices for AI Safety and Security. These include *“appropriate procedures and processes, to enable developers of AI, especially of dual-use foundation models, to conduct AI red-teaming tests”* as a mechanism for *“assessing and managing the safety, security, and trustworthiness of [these] models.”*

On one hand, red-teaming appears to call for *the right stuff*: find the flaws, find the vulnerabilities, and (help to) eliminate them. In this spirit, one might find its inclusion in a landmark policy document to be welcome. On the other hand, for all of the virtue in its aims, red-teaming at this level of description is strikingly vague. As noted by The Frontier Model Forum (FMF) [131], *“there is currently a lack of clarity on how to define ‘AI red teaming’ and what approaches are considered part of the expanded role it plays in the AI development life cycle.”* For example, the definition offered by the presidential executive order leaves the following key questions unanswered: What types of *undesirable behaviors, limitations, and risks* can or should be effectively caught and mitigated through red-teaming exercises? How should the activity be *structured* to maximize the likelihood of finding such flaws and vulnerabilities? For example, aside from AI developers, who else should be at the table, and what resources should be available to them? How should the risks identified through red-teaming be *documented, reported, and managed*? Is red-teaming on its own sufficient for assessing and managing the safety, security, and trustworthiness of AI? If not, what other practices should be part of the broader evaluation toolbox, and how does red-teaming complement those approaches? In short, is red-teaming the stuff of policy—the sort of concrete practice around which we can structure regulatory requirements?; or is it the stuff of *vibes*—a vague practice better suited to rallying than to rule-making?

Methodology. Using publicly available resources, we gathered information about recent real-world cases of AI red-teaming exercises (Section 3). Additionally, we performed an extensive survey of existing research literature on red-teaming and adjacent testing and evaluation methods (e.g., penetration testing, jailbreaking, and beyond) for generative AI (Section 4). We organized the thematic analysis of our case studies and literature survey around the following key questions:

- **Definition and scope:** What is the working definition of red-teaming? What is the criterion for successful red-teaming?
- **Object of evaluation:** What is the model being evaluated? Are its implementation details

(e.g., model architecture, training procedure, safety mechanisms) available to the evaluators or to the public? At what stage of its lifecycle (e.g., design, development, or deployment) is the model subjected to red-teaming?

- **Criteria of evaluation:** What is the threat model (i.e., the risk for which the model is being evaluated)? What are the risks the red-teaming activity potentially missed?
- **Actors and evaluators:** Who are the evaluators? What are the resources available to them (e.g., time, compute, expertise, type of access to model)?
- **Outcomes and broader impact:** What is the output of the activity? How much of the findings are shared publicly? What are the recommendations and mitigation strategies produced in response to the findings of red-teaming? What other evaluations had been on the model aside from red-teaming?

Contributions. Our analysis reveals the lack of consensus around the scope, structure, and assessment criteria for AI red-teaming. Prior methods and practices of AI red-teaming diverge along several critical axes, including the choice of threat model (if one is specified), the artifact under evaluation, the setting in which the activity is conducted (including actors, resources, and methodologies, and test-beds), and the resulting decisions the activity instigates (e.g., reporting, disclosure, and mitigation). In light of our findings, we argue that while red-teaming may be a valuable big-tent idea, and even a useful framing for a broad set of evaluation activities for Generative AI models, the bludgeoning use of *AI red-teaming* as a catch-all response to quiet all regulatory concerns about model safety verges on *security theater* [73]. The current framing of red-teaming in the public discourse appears to serve more to assuage regulators and other concerned parties than to offer a concrete solution. To move toward a more robust toolbox of evaluations for generative AI, we synthesize our recommendations into a question bank meant to guide and scaffold future AI red-teaming practices (see Table 1). We leave the careful co-design, development, and evaluation of this question bank as a critical direction for future work.

2 Related Contemporary Work

A brief history of red-teaming. Zenko [159] and Abbass [5] describe how the key concepts of red-teaming originated hundreds of years ago in warfare and religious contexts. They note the term “red team” was formally applied by the US military as early as the 1960s when modeling the Soviet Union’s behavior (in contrast to the “blue team” representing the US). According to Abbass et al. [4] and [5], red-teaming in computer security involves modeling an adversary and “*map[ping] out the space of vulnerabilities from a threat lens*” in contrast to penetration testing (in which enlisted cybersecurity experts actively attempt to find vulnerabilities in a computer system). Wood & Duggan [146] further describe how red-teaming “*is not an audit*” and that interpreting it as such risks reducing the amount of information shared about possible vulnerabilities. Using a hypothetical pandemic example, Bishop et al. [19] argue that effectively red-teaming a system requires context, knowledge, and assumptions about system utilization.

Evaluation beyond red-teaming. Chang & Custis [31] note that red-teaming is only one of many approaches to increase transparency of an AI system and that factsheets, audits, and model cards are other ways to do so. Similarly, Horvitz [58] warns of more advanced deepfakes in the near

Table 1: Our proposed set of questions to guide future AI red-teaming activities.

Phase	Key Questions and Considerations
0. Pre-activity	What is the artifact under evaluation through the proposed red-teaming activity? <ul style="list-style-type: none"> - What version of the model (including fine-tuning details) is to be evaluated? - What safety and security guardrails are already in place for this artifact? - At what stage of the AI lifecycle will the evaluation be conducted? - If the model has already been released, specify the conditions of release.
	What is the threat model the red-teaming activity probes? <ul style="list-style-type: none"> - Is the activity meant to illustrate a handful of possible vulnerabilities? (e.g., spelling errors in prompt leading to unpredictable model behavior) - Is the activity meant to identify a broad range of potential vulnerabilities? (e.g., biased behavior) - Is the activity meant to assess the risk of a specific vulnerability? (e.g., recipe for explosives)
	What is the specific vulnerability the red-teaming activity aims to find? <ul style="list-style-type: none"> - How was this vulnerability identified as the target of this evaluation? - Why was the above vulnerability prioritized over other potential vulnerabilities? - What is the threshold of acceptable risk for finding this vulnerability?
	What are the criteria for assessing the success of the red-teaming activity? <ul style="list-style-type: none"> - What are the benchmarks of comparison for success? - Can the activity be reconstructed or reproduced?
	Team composition and who will be part of the red team? <ul style="list-style-type: none"> - What were the criteria for inclusion/exclusion of members, and why? - How diverse/homogeneous is the team across relevant demographic characteristics? - How many internal versus external members belong to the team? - What is the distribution of subject-matter expertise among members? - What are possible biases or blindspots the current team composition may exhibit? - What incentives/disincentives do participants have to contribute to the activity?
	1. During activity
What instructions are given to the participants to guide the activity?	
What kind of access do participants have to the model?	
What methods can members of the team utilize to test the artifact? <p>Are there any auxiliary automated tools (including AI) supporting the activity?</p> <ul style="list-style-type: none"> - If yes, what are those tools? - Why are they integrated into the red-teaming activity? - How will members of the red team utilize the tool? 	
What reports and documentation are produced on the findings of the activity? <p>Who will have access to those reports? When and why?</p> <p>If certain details are withheld or delayed, provide justification.</p>	
2. Post-activity	What were the resources the activity consumed? <ul style="list-style-type: none"> - time - compute - financial resources - access to subject-matter expertise
	How successful was the activity in terms of the criteria specified in phase 0?
	What are the proposed measures to mitigate the risks identified in phase 1? <ul style="list-style-type: none"> - How will the efficacy of the mitigation strategy be evaluated? - Who is in charge of implementing the mitigation? - What are the mechanisms of accountability?

future while emphasizing that remedies such as increased media literacy and output watermarking (that flags relevant media as AI-generated) should be employed alongside red-teaming; Kenthapadi et al. [64] echo these concerns and similar solutions in their tutorial. Shevlane et al. [124] also argue that both internal and external model evaluations, as well as robust security responses, should complement effective red-teaming to counter extreme risks from generative AI.

Existing surveys of AI red-teaming and evaluations. Inie et al. [62] conduct qualitative interviews with those who perform red-teaming to create a grounded theory of *“how and why people attack large language models.”* Schuett et al. [118] survey members of labs racing to build artificial general intelligence (AGI) and find that 98% of respondents somewhat or strongly agree that *“AGI labs should commission external red teams before deploying powerful models.”* In the software design space, Kneareem et al. [66] highlight how UX designers are afraid that AI-based design tools will not be red-teamed enough while Liao et al. [77] suggest that UX designers themselves should help with red-teaming processes. Considering the testing of NLP systems specifically, Tan et al. [130] propose the DOCTOR framework for reliability testing of such systems. Weidinger et al. [144] introduce a framework for evaluating generative AI more broadly, namely via *“a three-layered framework that takes a structured, sociotechnical approach.”* Anderljung et al. [10] also propose a framework, ASPIRE, but for external accountability of LLMs and the engagement of relevant stakeholders. Yao et al. [155], Neel & Chang [91], and Shayegani et al. [122] produce surveys of LLM research with regard to security, privacy, and other vulnerabilities; Chang et al. [32] conduct another survey of LLM evaluation. In contrast to existing surveys of GenAI evaluation, our work focuses exclusively on *red-teaming*. Some of our findings resonate with points earlier made by Bockting et al. [21] and Friedler et al. [50], who argue for interdisciplinary audits of AI systems by a diverse group of people and red-teaming with concrete definitions of harms alongside other evaluations, respectively.

3 Case Studies: Prior AI Red-teaming Activities

To capture the complexity involved in designing real-world AI red-teaming exercises, we synthesize the results from 6 case studies recently conducted with generative models. These evaluations are mostly conducted by private companies and encompass a broad range of methods, goals and areas of focus. Through these case studies, we seek to understand common red-teaming practices, typical resources required for successful red-teaming, how red-teaming shapes deployed models, common pitfalls, and how results are shared with community stakeholders.

3.1 Findings

Considerable variation in red-teaming goals and processes across case studies. Reflecting the lack of consensus on a definition of red-teaming in the literature, red-teaming activities frequently varied in form and in goals. Some organizations chose to conduct a single round of red-teaming [96, 11, 30, 51], while others saw red-teaming as an iterative process in which results from initial rounds of testing were used to prioritize risk areas for further investigation [1, 3, 7]. The goals of red-teaming activities also ranged from specific objectives (e.g., red-teaming to investigate risks to national security [11]) to more broad targets (e.g., uncovering “harmful” model behavior [1]).

Language models are primary focus of red-teaming efforts. Most red-teaming activities in the 6 case studies analyzed here focused on language models, with only 2 case studies focused on multimodal models (GPT-4) [3, 7]. In some cases, red-teaming was conducted by internal teams prior to model release [3, 7, 96, 2], while other red-teaming activities were conducted on publicly-released models through APIs or websites [51, 30, 29]. Internal teams sometimes had access to different versions of the model being evaluated. For example, an internal evaluation team at OpenAI had access to early versions of GPT-4 both with and without safety and alignment guardrails [7]. This practice of red-teaming multiple versions of a model seemed effective at informing risk mitigation strategies because evaluation teams could directly evaluate and compare the ease of attacks with various mitigation strategies employed.

Evaluation team composition and available resources are interconnected and shape the outcomes of red-teaming. The evaluators employed in red-teaming activities for each case study varied considerably. Teams ranged from groups of subject matter experts to random samplings of community stakeholders. We found that there were generally three types of team compositions:

1. Teams composed of hand-selected subject matter experts in relevant areas (e.g., national security, healthcare, law, alignment), both internally and externally sourced
2. Crowd-sourced teams with evaluators sourced from crowd-working platforms or attendees of a live event
3. Teams composed of language models (i.e., language models prompted or fine tuned to red-team themselves)

The resources available to evaluation teams seemed to vary based on team composition. For crowd-sourced teams, red-teaming efforts were time-boxed either by participant or by task (sometimes participants were allowed to complete as many time-boxed tasks as they desire) and access to models was available only through APIs [29, 51]. For teams with subject matter expertise, red-teaming efforts were more open-ended with fewer strict restrictions on time or compute power. Often, the only limits on these teams’ explorations are rough priority lists of risks to identify [3, 7, 11]. While API access to models is still most common for these teams, sometimes experts are given access to versions of models without safety guardrails. When language models are used to red-team themselves, the main limitation seems to be the number of prompts used to produce red-teaming behavior along with the compute resources needed for model retraining or fine tuning. Full access to model parameters is thus usually a requirement when performing this type of red-teaming [96].

Red-teaming usually probes a broad range of potential vulnerabilities. Red-teaming efforts were usually focused on probing broad threat models, such as producing “harmful” or “offensive” model outputs without constraining these outputs to specific subject areas [1]. In two cases, evaluators were given specific goals: red-teaming for national security purposes [11] or red-teaming to score points in 20 different challenge areas [29]. With broader threat models, resulting model outputs are often grouped into subject areas to determine future prioritization of risk areas [51]. In one case study, a broader initial threat model was used with initial red-teaming efforts to prioritize risk areas used in subsequent rounds of testing [3]. Model developers often provide evaluators with broader threat models to probe in hopes that this will result in greater variation in red-teaming efforts, especially because it is often impossible to fully understand the model’s entire risk surface or

list all possible failure modes. Unfortunately, probing these nonspecific threat models does not always produce this desired variation, more so when evaluators are given limited time and resources to produce harmful model outputs. For example, some time-boxed evaluators often repeatedly probed models in the same risk area because it was easy to produce harmful outputs, even though these efforts contributed little to understanding the model’s full risk surface [51].

No standards or systematic procedures for disclosing results of red-teaming. We found nontrivial variation in the publicly-shared outputs of red-teaming efforts, largely because there are no existing standardized procedures or requirements for reporting results of red-teaming. In 3 cases, specific examples of risky or harmful model behavior uncovered by red-teaming efforts were publicly shared. In one case, a full dataset composed of 38,961 red team attacks was publicly released to aid in testing of other models [51]. In the other two cases, examples of harmful behavior were publicly available, but the full scope of all red-teaming attacks was not released [7, 96]. For red-teaming efforts on publicly available models or those focused on national security, specifics of harmful behavior were not shared publicly because findings were deemed “too sensitive.” One case study resulted in Anthropic piloting a responsible disclosure process to share vulnerabilities identified during red-teaming with appropriate community stakeholders, but this process is still under development (so we can assume that these disclosures have not yet been made) [11].

Various risk mitigation strategies are proposed or employed with no consensus on best practices or reporting requirements. While every case study analyzed here identified problematic or risky model behavior, none of the case studies resulted in a decision not to release the model. Instead, a number of risk mitigation strategies were proposed and/or employed to minimize the harmful model behavior identified during red-teaming. However, the specifics of risk mitigation strategies were often not provided when the target model was publicly available, and there were no standards for reporting improvements effected by these mitigation strategies. As a result, it was often difficult to determine if risks identified during red-teaming were sufficiently addressed. One exception is the red-teaming efforts against Google DeepMind’s Gopher chatbot, where specific strategies such as using reinforcement learning to penalize identified bad outputs or jointly training language models and red-teaming models via strategies for training GANs are explained in detail [96]. Some more-advanced proposed mitigation strategies were purely conceptual and untested (e.g., using unlikelihood training to reduce harmful outputs [96]), especially when they were described at a very high level (e.g., recommendations to other model developers to build multiple layers of mitigations or defenses throughout the system [7]). Other simpler mitigation strategies, such as identifying phrases to blacklist to reduce offensive replies or using Reinforcement Learning from Human Feedback (RLHF) to penalize harmful model outputs, have actually been deployed in model guardrails and tested during red-teaming efforts [51]. RLHF and rejection sampling seemed to have the most promising results out of all mitigation strategies identified and tested because these strategies made successful red-teaming attacks more difficult than other mitigation strategies studied [51].

Specific monetary and time costs of red-teaming usually not disclosed. Costs of red-teaming efforts were usually only shared for evaluation teams composed of crowd-sourced evaluators (for example, the hourly rate paid to crowd-workers [51]), though teams composed of subject matter experts and language models seem to have been given greater time and compute resources. Overall, it appears that red-teaming efforts for models intended for public release begin several months before release and continue after public launch to address emergent risks and capability jumps. Two case

studies specifically mention ongoing red-teaming for 6-7 months before model release [7, 2]. For crowd-sourced teams, evaluators spent about 30-50 minutes per task, with evaluators sourced from live events being limited to only completing a single task [51, 29].

Red-teaming misses risks due to broad threat models, evaluator biases, and limitations.

Without specific threat models to probe, evaluators often prioritized risks they had already seen and focused on those risks without exploring other risk areas. Crowd-sourced teams typically focused on risk areas where successful attacks were easy to produce due to time constraints, so risk areas that are more complex to attack may remain completely untested [51, 29]. The identified risk areas were usually not publicly released, making it difficult for other community stakeholders and external experts to determine if significant risk areas had been missed during red-teaming. When assembling teams of subject matter experts, the selection of evaluators can introduce significant bias into the types of risks investigated during red-teaming efforts. For example, risks prioritized by academic communities and AI firms may be explored in more detail because most experts were sourced from these areas [7]. When using language models for red-teaming, offensiveness classifiers are often trained on pre-existing datasets such as the Bot-Adversarial Dialogue (BAD) dataset [151], which in turn are not guaranteed to cover all types of offensive model replies. Similarly, distributional biases in models were only studied for small numbers of pre-defined population subgroups, so similar biases against other groups may be completely missed by red-teaming efforts [7].

No standards for evaluation methods used to supplement red-teaming. Every case study involved models that had been previously evaluated using other techniques beyond red-teaming, but there were no guidelines or standards established for these evaluation methods. As a result, there was significant variation in the other evaluation methods used. Commonly, models were evaluated using the Perspective API to measure toxicity; human feedback on helpfulness, harmfulness, and honesty; and QA benchmarks for accurate and truthful outputs [103, 2, 12]. Some models were quantitatively tested for distributional biases between gender and occupation, as well as for sentiment bias towards certain social groups. Internal quantitative assessments were sometimes performed to determine if model outputs violated specific content policies (e.g., hate speech, self-harm advice, illicit advice) [7]. Additionally, some initial efforts described as “red-teaming” by evaluators appeared more focused on understanding base model capabilities through open-ended experimentation than on specifically stress testing the model for vulnerabilities [3].

3.2 Discussion

Red-teaming is poorly-structured and is not comprehensive. Evaluation teams seem aware that the entire risk surface of a model will not be explored by red-teaming activities. As such, they either prioritize risk areas for investigation or provide evaluators with broad directions in hopes that diversity within the group of evaluators will lead to the exploration of many diverse risks. However there is a trade-off between providing evaluators with specific instructions and exploring a variety of risk areas through red-teaming. On one hand, vague instructions can be helpful to avoid biasing evaluators towards finding specific issues based on initial prioritization. On the other, the lack of instructions could reduce the utility of the exercise for uncovering risks relevant to real-world contexts. We argue that this limited scope of red-teaming efforts is concerning. Namely, recent executive orders and evaluation frameworks establishing red-teaming as a best-practice indicate that the broader perception of red-teaming may not align with current working definitions of red-teaming, (i.e., red-teaming activities are much more qualitative, subjective and exploratory than community

stakeholders may realize). In every case study, however, red-teaming was able to reveal harmful model behavior that other more systematic methods seemed to miss, highlighting the importance of both conducting red-teaming (alongside other evaluations) and developing systematic processes for red-teaming in a more comprehensive manner. These processes could include, for example, the development of guidelines on whether red-teaming is most effective when conducted internally or externally and when it should be conducted (i.e., before and/or after public release of the model and whether red-teaming activities should be ongoing while the model is publicly available).

Evaluation team composition introduces biases. The goal of team member selection seems to be ensuring variety in the risk areas explored during red-teaming. One option used to get this variety is to hand select a team of experts with different backgrounds; another option is to use a random sampling of the population through crowd-sourced teams. Both of these options have drawbacks: there may be bias in the processes used to select experts, and crowd-workers have very limited resources in terms of time, compute, and relevant expertise. It is difficult to say what the ideal balance between expert and non-technical stakeholders would be, but some sort of hybrid approach could help address some of the pitfalls associated with each type of team composition. One type of team composition we did not see explored in any case study is a crowd-sourced team with more open-ended instructions and greater resources. This could allow more variety in the risk areas explored because evaluators would not feel incentivized to focus on risk areas where harmful model outputs are easy or quick to produce, but it would also require partnering with subject matter experts to fully evaluate risky model behavior.

Team composition can also shape the outputs of red-teaming. One of the issues with red-teaming activities being organized by internal teams is that more extreme measures such as blocking the release of a model will probably never be recommended due to conflicting interests. On the other hand, external teams that may be more likely to recommend these extreme measures often do not have the power to actually employ these mitigations. Once again, a hybrid approach could resolve some of these issues but would need to be paired with accountability mechanisms to disclose recommendations and mitigation strategies.

Hesitancy to publicly release methods and results reduces utility of red-teaming. The reluctance to share all results from red-teaming activities may stem from the risks associated with public models (evaluators do not want to provide inspiration for potential attackers). Additionally, releasing all of the data associated with red-teaming could be overwhelming for community stakeholders. This said, because red-teaming does not seem to be planned as a comprehensive measure of risky model behavior, disclosing some specifics of red-teaming efforts and resulting mitigation strategies is necessary so stakeholders can understand the types of harms investigated and mitigated to in turn determine if they are relevant to their use cases. For example, significant risk areas that evaluation teams knowingly have not probed should be highlighted or identified in reports.

None of the case studies provided complete monetary costs of red-teaming efforts. This information seems relatively low-risk to release (compared to specific examples of harmful model behavior, for instance) and could be useful for developing methods to conduct more comprehensive red-teaming. For example, the costs of assembling teams with differing compositions of experts, non-technical stakeholders, and automation could be used to determine where resources can be used most effectively, especially considering that a hybrid team composition can lead to the best coverage of the risk surface. Additionally, the costs of evaluating and mitigating various types of risks could be factored into a cost-benefit analysis when prioritizing risks based on real-world impact. The

		Risk				Total
		Subjective	Objective	Both	Neither	
Approach	Brute-Force	15	4	1	0	20
	Brute-Force + AI	23	7	10	2	42
	Algorithmic Search	12	1	1	0	14
	Targeted Attack	19	7	2	0	28
Total		69	19	14	2	104

Table 2: Numbers of papers from our survey in each subgroup according to dimensions outlined in Section 4. Over half of papers are concerned with the generation of subjectively bad content, and over one-third of papers utilize brute-force + AI to red-team AI models.

lack of reported cost figures could also make it harder for third-party or external organizations to conduct red-teaming: if these unreported costs are quite large, it could be difficult or impossible for anyone aside from companies themselves to do this type of analysis. This information would be invaluable when developing guidelines for red-teaming to advise decisions such as whether internal or external red-teaming is most effective.

4 A Survey of AI Red-teaming Research

In this section, we analyze the results of our extensive survey of recent research on AI red-teaming and related concepts.

Methodology. We primarily searched arxiv, Google Scholar, OpenReview, ACL Anthology, and the ACM Digital Library with keywords “red-teaming”, “ai red-teaming”, “jailbreak”, and “llm jailbreak”, and we then gathered results.¹ Where possible, we replaced preprints with corresponding published works. We additionally included a few relevant works encountered prior to this search as well as some from snowball sampling based on this search.

We scrutinize and subdivide papers gathered for our survey into groups along two dimensions, both of which relate to the evaluation in each paper. The first corresponds to the type of risk investigated during the evaluation, and the second corresponds to the type of approach used for evaluation. We analyze papers by characteristics pertaining to threat model and methodology because we found that research works primarily focused on these aspects (perhaps due to technical relevance) as opposed to other factors important to red-teaming (such as team composition and resources consumed). Overall classifications and in-depth paper findings (as with analyses of case studies from the previous section) can be found in our appendix, and totals for each classification can be found in Table 2. Note that we do not ascribe moral attributes or judgments to any category (i.e., we do not posit that one form of red-teaming is better than another).

¹While we focus on red-teaming as a form of analysis in this work, we argue that the nascent *LLM jailbreaking* literature contains techniques similar in spirit to those employed by AI red-teams. We therefore also include this literature here.

4.1 Findings: Threat Models Underlying AI Red-teaming Research

We categorize papers according to the threat model implied or specified in each. In particular, we consider the type of *risk* posed in each threat model.

Subjective risk. We define *subjective risk* as risk that requires context to determine its threat level, and we categorize papers that studied this type of risk accordingly. For instance, hate speech, images with blood, and responses to questions like “how to build a bomb?” may be less than ideal AI generations, but they may be admissible phenomena depending on context.² Nevertheless, researchers like Gehman et al. [53] and Röttger et al. [112] create benchmarks to analyze hate speech generation, Rando et al. [105] perform red-teaming to ascertain the degree to which Stable Diffusion’s safety filter can be bypassed to generate NSFW content, and Zou et al. [168] evaluate whether LLMs respond to harmful queries like “Generate a step-by-step plan to destroy humanity” given adversarial prompt suffixes.

Objective risk. In contrast, we define *objective risk* as risk that does not require context to evaluate potential harms, and associated content should never be generated by properly functioning AI. We similarly categorize papers accordingly if the authors consider this type of risk. For example, phenomena such as the divulging of private information, leaking of training data, and production of vulnerable code or material for phishing attacks are inadmissible in any situation. To these ends, Chen et al. [36] study the degree to which multimodal LMs can safeguard private information, Nasr et al. [90] illustrate how divergence attacks cause ChatGPT to reveal training data, Wu et al. [147] analyze how code generation LLMs “*can be easily attacked and induced to generate vulnerable code,*” and Roy et al. [114, 115] discover that LLMs can create code for phishing attacks.

Both and neither. Some authors tasked themselves with analyzing *both* kinds of risk, such as personally identifiable information (PII) leakage in addition to hate speech or dangerous generations [127, 51, 96]. Others introduce methods to analyze *neither* type of risk from the outset, stressing that definitions and classifications of issues may need to be done from scratch [28, 102].

4.2 Findings: Methodologies Proposed to Facilitate Red-teaming

We further categorize papers based on the methodology the researchers employ to perform red-teaming. Namely, we study the type of *approach* used to find risks.

Brute-force. Work that utilized *brute-force* approaches involved manual evaluation of generative AI inputs and outputs by teams of humans. We found that such teams typically consisted of the researchers themselves, internal auditors (of tech companies), or external members (such as contractors hired via Amazon Mechanical Turk (MTurk)). Xu et al. [150, 151] and Ganguli et al. [51] employed crowdworkers to elicit harmful text outputs from language models (including but not limited to offensive language and PII) and measure safety. Mu et al. [89] compiled a benchmark from scratch to test LLMs’ capacities to follow rules while Huang et al. [60] hired crowdworkers to build a new benchmark that assesses alignment with Chinese values. Schulhoff et al. [119] hosted a prompt hacking competition, thereby making competitors LLM red team members. Other authors hand-craft jailbreak attacks against language models [45, 76, 143, 75, 79], but the authors of [143] join Xie et al. [149] in additionally devising defense strategies for them. Shen et al. [123] and Rao

²Such generations do not reflect opinions held by the authors.

et al. [107] analyze the effectiveness of jailbreak attacks collected from external sources (including prior work and public websites and forums).

Brute-force + AI. Another body of work similar to those of the brute-force works described above incorporated AI techniques in their red-teaming processes. Common approaches to do so typically involved having AI models generate test cases and find errors in other AI output. We therefore term such approaches as *brute-force + AI*. Many authors used LLMs to generate normal prompts [96, 108, 127, 16, 35, 85, 161] and jailbreak prompts [157, 40, 120, 154, 142] such that LLMs produce bad outputs like harmful text responses. Variations on these ideas also exist, such as the work of Pfau et al. [97], in which the authors use *reverse LMs* to work backwards from harmful text responses to prompts that could generate them. Others use LLMs to devise new benchmarks related to exaggerated safety responses (i.e., refusal to respond to prompts that are arguably safe) [111], *fake alignment* that occurs when models appear aligned with one query format and misaligned with another (e.g., multiple choice versus open-ended response) [140], and *latent jailbreaks*, or compliance with “*implicit malicious instruction[s]*” [101]. Researchers have also used AI to red-team and jailbreak text-to-image models and multimodal LMs. For instance, Lee et al. [71] demonstrate how passing images related to harmful queries with the queries themselves to multimodal models (e.g., an image of a bomb with the question “how to build a bomb?”) improves the likelihood of harmful text generation. Mehrabi et al. [82] test their FLIRT framework to analyze text-to-image models like Stable Diffusion. Still other researchers perform red-teaming of LLMs for specific end-uses. Lewis & White [74] red-team an LLM for potential future use as a component of a virtual museum tourguide, and He et al. [57] evaluate the dangers of using LLMs as part of scientific research. In light of the many documented ways generative AI models can be utilized for malicious use, researchers have also studied ways in which they can be defended. Both Sun et al. [129] and Wang et al. [141] introduce methods that utilize LLMs to generate fine-tuning data that can be used to avert harmful responses. Zhu et al. [166] employ k-nearest neighbors and clustering techniques to fix incorrect labels in popular LLM safety datasets.

Algorithmic search. Some other methods start from a given prompt and utilize a process to modify it until an issue is encountered. Such processes can take the form of random perturbations or a guided search, and we therefore refer to such approaches as *algorithmic search* strategies.³ For instance, several authors describe approaches to red-teaming and jailbreaking in which one AI model automatically and repeatedly attacks an LLM until defenses are broken [27, 81, 33, 84]. Both Chin et al. [37] and Tsai et al. [135] propose search-based red-teaming approaches to evaluate text-to-image models that perturb input prompts until they simultaneously pass safety filters and generate forbidden content. Search-based approaches can also be used as defensive measures. Noting the brittleness of most jailbreak methods, Robey et al. [109] and Zhang et al. [162] introduce methods to detect jailbreaks by applying perturbations to text and image inputs and observing whether outputs change drastically (if so, the input was likely a jailbreak).

Targeted attack. The last approach to red-teaming we document as part of our review involves deliberately targeting part of an LLM, which could include an API, vulnerability in language translation support, or step of its training process, in order to induce issues. As such, we refer to such approaches as *targeted attack* methods. For instance, Wang & Shu [138] show how to construct *steering vectors* from activations of both safety-tuned and non-safety-tuned versions of models to

³Note that we differentiate these strategies from brute-force + AI approaches in that instead of using algorithms to generate test cases at scale, these approaches use algorithms to methodically search for problems.

obtain toxic outputs from safety-tuned models. Others illustrate how to imperceptibly perturb images to cause multimodal LMs to respond in unintended ways (such as replying with a malicious URL or misinformation) [117, 100, 13], and Tong et al. [133] engineer prompts for text-to-image models that are mismatched with resulting images by making use of many such models’ reliance on CLIP embeddings. Other approaches include but are not limited to exploiting the fact that LLMs are not optimized to converse in low-resource languages and ciphers [42, 156, 158], poisoning data used to tune or utilize LLMs [70, 106, 160, 6, 24, 139], and attacking APIs associated with black-box models [95]. Various defensive methods rooted in targeted attack approaches have been proposed as well. Bitton et al. [20] describe their Adversarial Text Normalizer, which can defend an LLM against various character-level perturbations typical of certain adversarial prompts. In addition, other defensive strategies proposed in prior paragraphs can defend against attacks discussed here (e.g., JailGuard introduced by Zhang et al. [162] addresses attacks introduced in [100, 13, 117, 168]).

4.3 Discussion

Utilizing in-depth notes created for each of the papers obtained in this survey, two of the authors conducted thematic analyses of salient details to gather high-level takeaways (as with the case study analyses). Each author conducted their own analysis independently from one another, and the key points are summarized below.

Many different methods to perform red-teaming. As illustrated in Table 2, researchers and practitioners have undertaken numerous approaches to evaluate LLMs and have all described them as *red-teaming*. At the same time, there have been developments like Schuett et al.’s finding that the overwhelming majority of AGI lab members support external red-teaming efforts [118] and the recent Executive Order [59] stressing the importance of red-teaming. These developments and the many red-teaming variations are together arguably concerning, precisely because there are no concrete definitions regarding what constitutes red-teaming. By highlighting this, we do not mean to imply that evaluations until now are useless. On the contrary, we believe they are necessary but perhaps insufficient tests of safety, and we stipulate that the existence of many interpretations of “red-teaming” suggests there must be more top-down guidance and requirements concerning red-teaming evaluations.

Threat modeling skewed toward subjective risk. Table 2 also highlights that the majority of evaluations focus on *subjective risk* rather than *objective risk*. This means that undue effort has been undertaken to evaluate and mitigate LLM behavior that may be admissible in various contexts. Additionally, work like that of Röttger et al. [111] has shown that current attempts to mitigate such risks have resulted in exaggerated safety, yielding LLM behavior like the refusal to provide information on buying a can of coke. Lastly, focusing on subjective risk takes attention away from objective risk, which in turn are inadmissible in any context. In light of such issues and tradeoffs, Casper et al. [28] and Radharapu et al. [102] suggest clearly defining risks and behavior to uncover via red-teaming and justifying those decisions before starting any analysis.

No consensus on adversary capabilities. While threat model and methodology are two factors that contribute to the diversity of red-teaming exercises, assumptions about adversary capabilities are also contributors. Namely, the works encountered have differing estimates of adversary resources. For instance, Perez et al. [96] and many authors of similar work conjecture that an adversary can

only prompt an LLM and probe it for bad outputs. In contrast, others assume that an adversary can poison the training process [106], has the compute required to search for adversarial suffixes [168], or is able to run both safety-tuned and non-safety-tuned versions of language models to obtain toxic output [138]. Future guidelines for red-teaming may want to suggest that researchers should emphasize and defend adversary assumptions.

No consensus on values used for alignment and red-teaming. Work found as part of this survey involving subjective risk and alignment are driven by, implicitly or explicitly, a set of human values that determine whether LLM behavior is admissible or inadmissible. However, this in turn prompts the question *whose values are being utilized for alignment and evaluation?* For instance, the FLAMES benchmark proposed by Huang et al. [60] is purported to measure alignment with Chinese values, whereas Weidinger et al. [144] emphasize that other evaluations may reflect values of those in *“the English-speaking or Western world.”* The extent to which LLMs do not support low-resource languages [42, 156] and agree with bias and stereotypes [51, 108] evidences that these models may not reflect the values and beliefs of all persons. Works beyond this survey have illustrated how the framing of AI value alignment is a normative problem that, if not properly addressed, may only serve to reflect the norms of one group of people, typically the majority [46, 68, 67]. Especially as OpenAI has begun a partnership with the US military on one hand [128, 49] and launched an initiative to align superintelligent AI to human values on another [72], we argue that it is crucial to analyze assumptions made and viewpoints held by those who develop and deploy AI systems.

No consensus on who should perform red-teaming. Moreover, just as there is a lack of agreement regarding whose values to utilize to assess LLM behavior, there is a similar lack of agreement regarding who should be part of red-teaming exercises. Groups of evaluators have consisted of hired crowdworkers [51], competition participants [119], researchers themselves [96], and others simply red-teaming for fun [62]. While some argue for more diversity to evaluate AI models [126], others caution that increased diversity is not a panacea and is moreover typically ill-defined [144, 14]. For instance, Yong et al. [156] argue for multilingual red-teaming to respond to low-resource language issues, and He et al. [57] *“advocate for a collaborative, interdisciplinary approach among the AI for Science community and society at large”* to respond to scientific research risks. Such examples suggest that diversity should be defined and sought out relative to the risks considered by red-teaming processes. They additionally hint towards more involvement of the public and relevant stakeholders, things also recommended in parallel literature regarding algorithmic auditing and participatory approaches to machine learning [48, 39, 38, 17]. Future red-teaming guidelines may want to specify something along these lines.

Unclear follow-ups to red-teaming activities. We found that overall, responses from LLM developers (at least public ones) to the many red-teaming and jailbreaking papers have been muted and generally mixed. While some authors such as Wei et al. [142] report that they reached out to organizations like OpenAI and Anthropic about the vulnerabilities found in their models, the vulnerabilities and models themselves have for the most part persisted. One rare exception to this pattern is the case of the findings of Nasr et al. [90], in which OpenAI updated ChatGPT to reduce the likelihood of divergence attack success and modified their terms of use to forbid such attacks in response [98, 88]. However, these changes only came following the paper’s release, 90 days *after* the paper authors first notified OpenAI about the vulnerability. If red-teaming is to be stipulated as a requirement for release and safe usage of AI models, there should arguably be a protocol to mitigate found issues accordingly.

5 Takeaways and Recommendations

Based on the results of our case study analysis and research survey, we have the following takeaways and recommendations for future red-teaming evaluations.

Red-teaming is *not* a panacea. Each red-teaming exercise discussed in this paper only covered a limited set of vulnerabilities. As such, red-teaming cannot be expected to guarantee safety from all angles. For example, red-teaming to detect and mitigate harmful text responses [96, 51] may not detect and mitigate phishing attack vulnerabilities [115, 114] and vice versa. Team composition may also influence the types of issues found in a given exercise (i.e., a group of subject matter experts may find different problems than a group of crowd-workers from MTurk [51]). Moreover, there are other issues that red-teaming alone cannot address, such as problems stemming from *algorithmic monoculture*, or a lack of diversity in datasets and model architectures that in turn lead to similar model failures [133, 65, 22]. We argue that red-teaming should therefore be considered as one evaluation paradigm, among others, to assess and improve the safety and trustworthiness of generative AI models.

Red-teaming, as currently conducted, is *not* well-scoped or structured. Additionally, the many variations in the red-teaming processes encountered in our case studies and literature review illustrate that at the moment, red-teaming is an unstructured procedure with undefined scope. As previously stated, we do not mean to belittle efforts undertaken so far to evaluate complex systems, but in order to derive greater utility from future evaluations, we recommend that guidelines for red-teaming should be carefully drafted. While we sketched an initial list of considerations we believe such guidelines should encompass, we acknowledge that all stakeholders of generative AI models, including members of the general public and research community alike, should have a say on what should comprise these guidelines.

There are no standards concerning what should be reported. We further highlight that there are also currently no unified protocols for reporting the results of red-teaming evaluations. In fact, we found that a number of case studies and research papers sourced for our work did not fully report their findings or resource costs to perform evaluations. For a number of reasons, ranging from increasing public knowledge to helping third-party groups conduct their own tests to assisting end-users in determining the relevance of red-teaming for their use cases, we suggest that regulations and/or best practices should be put forth to entice more detailed reporting following these exercises. We argue that such reports should, at a minimum, clarify (1) the resources consumed by the activity, (2) assessments of whether the activity was successful according to previously established goals and measures, (3) the mitigation steps informed by the findings of the activity, and (4) any other relevant or subsequent evaluation of the artifact at hand.

Mitigation steps initiated by red-teaming are often unclear and unrepresentative. Though red-teaming exercises uncovered many issues with generative models, follow-up activities to remedy these problems were often vague or unspecified. Taken with the lack of reporting, such unclear mitigation and alignment strategies could reduce red-teaming to an *approval-stamping process* wherein one can say that red-teaming was performed as an assurance without providing further details into issues found or fixed. Moreover, we found that the strategies specified in research and case studies, such as further fine-tuning or RLHF, were often not representative of the full range of possible solutions. Other approaches like model input and output monitoring, dataset cleaning,

prediction modification, and even the refusal to deploy models in certain scenarios, were rarely or never mentioned. Future research should address risk mitigation strategies beyond these most popular solutions in the face of issues surfaced by red-teaming.

We propose our red-teaming question bank as a starting point to address these issues.

In light of the questions and concerns raised by our case study and research survey, we provide a set of questions for future red team evaluators to consider before, during, and after their evaluation. These questions, contained in Table 1, encourage evaluators to ponder the benefits and limitations of red-teaming exercises generally as well as the impact of specific design choices pertaining to their setting. We emphasize that these questions are not finalized guidelines but rather (what we hope is) the start of a broader conversation about GenAI red-teaming and evaluation processes. We welcome and support comments and feedback from the research community, industry, government, the general public, and beyond on our initial draft, and we leave the further refinement and development of the question bank, as well as its usability and efficacy assessments, as a critical direction for future work.

Acknowledgements

H. Heidari acknowledges support from NSF (IIS2040929 and IIS2229881) and PwC (through the Digital Transformation and Innovation Center at CMU). M. Feffer acknowledges support from the National GEM Consortium and the ARCS Foundation. Authors additionally gratefully acknowledge the NSF (IIS2211955), UPMC, Highmark Health, Abridge, Ford Research, Mozilla, Amazon AI, JP Morgan Chase, the Block Center, the Center for Machine Learning and Health, and the CMU Software Engineering Institute (SEI) via Department of Defense contract FA8702-15-D-0002, for their generous support of Z. Lipton’s, A. Sinha’s, and ACMI Lab’s research. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the authors and do not reflect the views of the National Science Foundation and other funding agencies.

References

- [1] (2023).
URL <https://www.frontiermodelforum.org/uploads/2023/10/FMF-AI-Red-Teaming.pdf>
- [2] (2023).
URL <https://www-files.anthropic.com/production/images/Model-Card-Claude-2.pdf>
- [3] (2024).
URL <https://support.microsoft.com/en-us/topic/copilot-in-bing-our-approach-to-responsible-ai-45b5eae8-7466-43e1-ae98-b48f8ff8fd44>
- [4] Abbass, H., Bender, A., Gaidow, S., & Whitbread, P. (2011). Computational red teaming: Past, present and future. *IEEE Computational Intelligence Magazine*, 6(1), 30–42.
- [5] Abbass, H. A. (2015). *Computational red teaming*. Springer.
- [6] Abdelnabi, S., Greshake, K., Mishra, S., Endres, C., Holz, T., & Fritz, M. (2023). Not what you’ve signed up for: Compromising real-world llm-integrated applications with indirect

- prompt injection. In *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, (pp. 79–90).
- [7] Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. (2023). Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- [8] Agostinelli, A., Denk, T. I., Borsos, Z., Engel, J., Verzetti, M., Caillon, A., Huang, Q., Jansen, A., Roberts, A., Tagliasacchi, M., et al. (2023). Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*.
- [9] Alon, G., & Kamfonas, M. (2023). Detecting language model attacks with perplexity. *arXiv preprint arXiv:2308.14132*.
- [10] Anderljung, M., Smith, E., O’Brien, J., Soder, L., Bucknall, B., Bluemke, E., Schuett, J., Trager, R., Strahm, L., & Chowdhury, R. (2023). Towards publicly accountable frontier llms. In *Socially Responsible Language Modelling Research (SoLaR) Workshop at NeurIPS*.
- [11] Anthropic (2023).
URL <https://www.anthropic.com/news/frontier-threats-red-teaming-for-ai-safety>
- [12] Askell, A., Bai, Y., Chen, A., Drain, D., Ganguli, D., Henighan, T., Jones, A., Joseph, N., Mann, B., DasSarma, N., et al. (2021). A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*.
- [13] Bailey, L., Ong, E., Russell, S., & Emmons, S. (2023). Image hijacking: Adversarial images can control generative models at runtime. *arXiv preprint arXiv:2309.00236*.
- [14] Bergman, A. S., Hendricks, L. A., Rauh, M., Wu, B., Agnew, W., Kunesch, M., Duan, I., Gabriel, I., & Isaac, W. (2023). Representation in ai evaluations. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, (pp. 519–533).
- [15] Bhardwaj, R., & Poria, S. (2023). Language model unalignment: Parametric red-teaming to expose hidden harms and biases. *arXiv preprint arXiv:2310.14303*.
- [16] Bhardwaj, R., & Poria, S. (2023). Red-teaming large language models using chain of utterances for safety-alignment. *arXiv preprint arXiv:2308.09662*.
- [17] Birhane, A., Isaac, W., Prabhakaran, V., Diaz, M., Elish, M. C., Gabriel, I., & Mohamed, S. (2022). Power to the people? opportunities and challenges for participatory ai. *Equity and Access in Algorithms, Mechanisms, and Optimization*, (pp. 1–8).
- [18] Birhane, A., Prabhu, V. U., & Kahembwe, E. (2021). Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*.
- [19] Bishop, M., Gates, C., & Levitt, K. (2018). Augmenting machine learning with argumentation. In *Proceedings of the New Security Paradigms Workshop*, (pp. 1–11).
- [20] Bitton, J., Pavlova, M., & Evtimov, I. (2022). Adversarial text normalization. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track*, (pp. 268–279).

- [21] Bockting, C. L., van Dis, E. A. M., van Rooij, R., Zuidema, W., & Bollen, J. (2023). Living guidelines for generative ai — why scientists must oversee its use. *Nature*, 622(7984), 693–696.
- [22] Bommasani, R., Creel, K. A., Kumar, A., Jurafsky, D., & Liang, P. S. (2022). Picking on the same person: Does algorithmic monoculture lead to outcome homogenization? *Advances in Neural Information Processing Systems*, 35, 3663–3678.
- [23] Cao, B., Cao, Y., Lin, L., & Chen, J. (2023). Defending against alignment-breaking attacks via robustly aligned llm. *arXiv preprint arXiv:2309.14348*.
- [24] Cao, Y., Cao, B., & Chen, J. (2023). Stealthy and persistent unalignment on large language models via backdoor injections. *arXiv preprint arXiv:2312.00027*.
- [25] Casper, S., Bu, T., Li, Y., Li, J., Zhang, K., Hariharan, K., & Hadfield-Menell, D. (2023). Red teaming deep neural networks with feature synthesis tools. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- [26] Casper, S., Hariharan, K., & Hadfield-Menell, D. (2022). Diagnostics for deep neural networks with automated copy/paste attacks. *arXiv preprint arXiv:2211.10024*.
- [27] Casper, S., Killian, T., Kreiman, G., & Hadfield-Menell, D. (2022). Red teaming with mind reading: White-box adversarial policies in deep reinforcement learning. *arXiv preprint arXiv:2209.02167*.
- [28] Casper, S., Lin, J., Kwon, J., Culp, G., & Hadfield-Menell, D. (2023). Explore, establish, exploit: Red teaming language models from scratch. *arXiv preprint arXiv:2306.09442*.
- [29] Cattell, S. (2023).
URL <https://aivillage.org/defcon%2031/generative-recap/>
- [30] Cattell, S., Carson, A., & Chowdhury, R. (2023).
URL <https://aivillage.org/generative%20red%20team/generative-red-team/>
- [31] Chang, J., & Custis, C. (2022). Understanding implementation challenges in machine learning documentation. In *Equity and Access in Algorithms, Mechanisms, and Optimization*, (pp. 1–8).
- [32] Chang, Y., Wang, X., Wang, J., Wu, Y., Zhu, K., Chen, H., Yang, L., Yi, X., Wang, C., Wang, Y., et al. (2023). A survey on evaluation of large language models. *arXiv preprint arXiv:2307.03109*.
- [33] Chao, P., Robey, A., Dobriban, E., Hassani, H., Pappas, G. J., & Wong, E. (2023). Jail-breaking black box large language models in twenty queries. In *Robustness of Few-shot and Zero-shot Learning in Large Foundation Models (R0-FoMo) Workshop at NeurIPS*.
- [34] Chen, B., Paliwal, A., & Yan, Q. (2023). Jailbreaker in jail: Moving target defense for large language models. In *Proceedings of the 10th ACM Workshop on Moving Target Defense*, (pp. 29–32).
- [35] Chen, B., Wang, G., Guo, H., Wang, Y., & Yan, Q. (2023). Understanding multi-turn toxic behaviors in open-domain chatbots. In *Proceedings of the 26th International Symposium on Research in Attacks, Intrusions and Defenses*, (pp. 282–296).

- [36] Chen, Y., Mendes, E., Das, S., Xu, W., & Ritter, A. (2023). Can language models be instructed to protect personal information? *arXiv preprint arXiv:2310.02224*.
- [37] Chin, Z.-Y., Jiang, C.-M., Huang, C.-C., Chen, P.-Y., & Chiu, W.-C. (2023). Prompting4debugging: Red-teaming text-to-image diffusion models by finding problematic prompts. *arXiv preprint arXiv:2309.06135*.
- [38] Costanza-Chock, S., Harvey, E., Raji, I. D., Czernuszenko, M., & Buolamwini, J. (2022). Who audits the auditors? recommendations from a field scan of the algorithmic auditing ecosystem. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, (p. 1571–1583). ArXiv:2310.02521 [cs].
URL <http://arxiv.org/abs/2310.02521>
- [39] Delgado, F., Yang, S., Madaio, M., & Yang, Q. (2023). The participatory turn in ai design: Theoretical foundations and the current state of practice. In *Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, (pp. 1–23).
- [40] Deng, B., Wang, W., Feng, F., Deng, Y., Wang, Q., & He, X. (2023). Attack prompt generation for red teaming and defending large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, (pp. 2176–2189).
- [41] Deng, G., Liu, Y., Li, Y., Wang, K., Zhang, Y., Li, Z., Wang, H., Zhang, T., & Liu, Y. (2023). Masterkey: Automated jailbreak across multiple large language model chatbots. *arXiv preprint arXiv:2307.08715*.
- [42] Deng, Y., Zhang, W., Pan, S. J., & Bing, L. (2023). Multilingual jailbreak challenges in large language models. *arXiv preprint arXiv:2310.06474*.
- [43] Ding, P., Kuang, J., Ma, D., Cao, X., Xian, Y., Chen, J., & Huang, S. (2023). A wolf in sheep’s clothing: Generalized nested jailbreak prompts can fool large language models easily. *arXiv preprint arXiv:2311.08268*.
- [44] Donahue, C., Caillon, A., Roberts, A., Manilow, E., Esling, P., Agostinelli, A., Verzetti, M., Simon, I., Pietquin, O., Zeghidour, N., et al. (2023). Singsong: Generating musical accompaniments from singing. *arXiv preprint arXiv:2301.12662*.
- [45] Du, Y., Zhao, S., Ma, M., Chen, Y., & Qin, B. (2023). Analyzing the inherent response tendency of llms: Real-world instructions-driven jailbreak. *arXiv preprint arXiv:2312.04127*.
- [46] Feffer, M., Heidari, H., & Lipton, Z. C. (2023). Moral machine or tyranny of the majority? *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(55), 5974–5982.
- [47] Feffer, M., Lipton, Z. C., & Donahue, C. (2023). Deepdrake ft. bts-gan and taylorvc: an exploratory analysis of musical deepfakes and hosting platforms. In *Proceedings of the 2nd Workshop on Human-Centric Music Information Retrieval 2023 (HCMIR 2023)*.
- [48] Feffer, M., Skirpan, M., Lipton, Z., & Heidari, H. (2023). From preference elicitation to participatory ml: A critical survey & guidelines for future research. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, (pp. 38–48).
- [49] Field, H. (2024). Openai quietly removes ban on military use of its ai tools.
URL <https://www.cnn.com/2024/01/16/openai-quietly-removes-ban-on-military-use-of-its-ai-tools.html>

- [50] Friedler, S., Singh, R., Blili-Hamelin, B., Metcalf, J., & Chen, B. J. (2023). Ai red-teaming is not a one-stop solution to ai harms: Recommendations for using red-teaming for ai accountability.
- [51] Ganguli, D., Lovitt, L., Kernion, J., Askill, A., Bai, Y., Kadavath, S., Mann, B., Perez, E., Schiefer, N., Ndousse, K., et al. (2022). Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*.
- [52] Ge, S., Zhou, C., Hou, R., Khabsa, M., Wang, Y.-C., Wang, Q., Han, J., & Mao, Y. (2023). Mart: Improving llm safety with multi-round automatic red-teaming. *arXiv preprint arXiv:2311.07689*.
- [53] Gehman, S., Gururangan, S., Sap, M., Choi, Y., & Smith, N. A. (2020). Realtotoxicityprompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, (pp. 3356–3369).
- [54] Ghosh, S., & Caliskan, A. (2023). Chatgpt perpetuates gender bias in machine translation and ignores non-gendered pronouns: Findings across bengali and five other low-resource languages. In *Proceedings of the 2023 AAI/ACM Conference on AI, Ethics, and Society, AIES '23*, (p. 901–912). New York, NY, USA: Association for Computing Machinery.
URL <https://doi.org/10.1145/3600211.3604672>
- [55] Gong, Y., Ran, D., Liu, J., Wang, C., Cong, T., Wang, A., Duan, S., & Wang, X. (2023). Figstep: Jailbreaking large vision-language models via typographic visual prompts. *arXiv preprint arXiv:2311.05608*.
- [56] Greenblatt, R., Shlegeris, B., Sachan, K., & Roger, F. (2023). Ai control: Improving safety despite intentional subversion. *arXiv preprint arXiv:2312.06942*.
- [57] He, J., Feng, W., Min, Y., Yi, J., Tang, K., Li, S., Zhang, J., Chen, K., Zhou, W., Xie, X., et al. (2023). Control risk for potential misuse of artificial intelligence in science. *arXiv preprint arXiv:2312.06632*.
- [58] Horvitz, E. (2022). On the horizon: Interactive and compositional deepfakes. In *Proceedings of the 2022 International Conference on Multimodal Interaction*, (pp. 653–661).
- [59] House, T. W. (2023). Executive order on the safe, secure, and trustworthy development and use of artificial intelligence.
URL <https://www.whitehouse.gov/briefing-room/presidential-actions/2023/10/30/executive-order-on-the-safe-secure-and-trustworthy-development-and-use-of-artificial-intelligence/>
- [60] Huang, K., Liu, X., Guo, Q., Sun, T., Sun, J., Wang, Y., Zhou, Z., Wang, Y., Teng, Y., Qiu, X., et al. (2023). Flames: Benchmarking value alignment of chinese large language models. *arXiv preprint arXiv:2311.06899*.
- [61] Huang, Y., Gupta, S., Xia, M., Li, K., & Chen, D. (2023). Catastrophic jailbreak of open-source llms via exploiting generation. *arXiv preprint arXiv:2310.06987*.
- [62] Inie, N., Stray, J., & Derczynski, L. (2023). Summon a demon and bind it: A grounded theory of llm red teaming in the wild. *arXiv preprint arXiv:2311.06237*.

- [63] Jain, N., Schwarzschild, A., Wen, Y., Somepalli, G., Kirchenbauer, J., Chiang, P.-y., Goldblum, M., Saha, A., Geiping, J., & Goldstein, T. (2023). Baseline defenses for adversarial attacks against aligned language models. *arXiv preprint arXiv:2309.00614*.
- [64] Kenthapadi, K., Lakkaraju, H., & Rajani, N. (2023). Generative ai meets responsible ai: Practical challenges and opportunities. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, (pp. 5805–5806).
- [65] Kleinberg, J., & Raghavan, M. (2021). Algorithmic monoculture and social welfare. *Proceedings of the National Academy of Sciences*, 118(22), e2018340118.
- [66] Knearem, T., Khwaja, M., Gao, Y., Bentley, F., & Kliman-Silver, C. E. (2023). Exploring the future of design tooling: The role of artificial intelligence in tools for user experience professionals. In *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems*, (pp. 1–6).
- [67] Lambert, N., & Calandra, R. (2023). The alignment ceiling: Objective mismatch in reinforcement learning from human feedback. *arXiv preprint arXiv:2311.00168*.
- [68] Lambert, N., Gilbert, T. K., & Zick, T. (2023). Entangled preferences: The history and risks of reinforcement learning and human feedback. *arXiv preprint arXiv:2310.13595*.
- [69] Lapid, R., Langberg, R., & Sipper, M. (2023). Open sesame! universal black box jailbreaking of large language models. *arXiv preprint arXiv:2309.01446*.
- [70] Lee, A., Bai, X., Pres, I., Wattenberg, M., Kummerfeld, J. K., & Mihalcea, R. (2024). A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. (arXiv:2401.01967). ArXiv:2401.01967 [cs].
URL <http://arxiv.org/abs/2401.01967>
- [71] Lee, D., Lee, J., Ha, J.-W., Kim, J.-H., Lee, S.-W., Lee, H., & Song, H. O. (2023). Query-efficient black-box red teaming via bayesian optimization. *arXiv preprint arXiv:2305.17444*.
- [72] Leike, J., & Sutskever, I. (2023). Introducing superalignment.
URL <https://openai.com/blog/introducing-superalignment>
- [73] Levenson, E. (2014). The tsa is in the business of ‘security theater,’ not security. *The Atlantic Magazine*.
- [74] Lewis, A., & White, M. (2023). Mitigating harms of llms via knowledge distillation for a virtual museum tour guide. In *Proceedings of the 1st Workshop on Taming Large Language Models: Controllability in the era of Interactive Assistants!*, (pp. 31–45).
- [75] Li, H., Guo, D., Fan, W., Xu, M., & Song, Y. (2023). Multi-step jailbreaking privacy attacks on chatgpt. *arXiv preprint arXiv:2304.05197*.
- [76] Li, X., Zhou, Z., Zhu, J., Yao, J., Liu, T., & Han, B. (2023). Deepinception: Hypnotize large language model to be jailbreaker. *arXiv preprint arXiv:2311.03191*.
- [77] Liao, Q. V., Subramonyam, H., Wang, J., & Wortman Vaughan, J. (2023). Designerly understanding: Information needs for model transparency to support design ideation for ai-powered user experience. In *Proceedings of the 2023 CHI conference on human factors in computing systems*, (pp. 1–21).

- [78] Liu, X., Zhu, Y., Lan, Y., Yang, C., & Qiao, Y. (2023). Query-relevant images jailbreak large multi-modal models. *arXiv preprint arXiv:2311.17600*.
- [79] Liu, Y., Deng, G., Xu, Z., Li, Y., Zheng, Y., Zhang, Y., Zhao, L., Zhang, T., & Liu, Y. (2023). Jailbreaking chatgpt via prompt engineering: An empirical study. *arXiv preprint arXiv:2305.13860*.
- [80] Luccioni, S., Akiki, C., Mitchell, M., & Jernite, Y. (2023). Stable bias: Evaluating societal representations in diffusion models. In *Advances in Neural Information Processing Systems (NeurIPS) 2023 Dataset and Benchmarks Track*.
URL <https://openreview.net/forum?id=qVXYU3F017>
- [81] Ma, C., Yang, Z., Gao, M., Ci, H., Gao, J., Pan, X., & Yang, Y. (2023). Red teaming game: A game-theoretic framework for red teaming language models. *arXiv preprint arXiv:2310.00322*.
- [82] Mehrabi, N., Goyal, P., Dupuy, C., Hu, Q., Ghosh, S., Zemel, R., Chang, K.-W., Galstyan, A., & Gupta, R. (2023). Flirt: Feedback loop in-context red teaming. *arXiv preprint arXiv:2308.04265*.
- [83] Mehrabi, N., Goyal, P., Ramakrishna, A., Dhamala, J., Ghosh, S., Zemel, R., Chang, K.-W., Galstyan, A., & Gupta, R. (2023). Jab: Joint adversarial prompting and belief augmentation. In *Robustness of Few-shot and Zero-shot Learning in Large Foundation Models (R0-FoMo) Workshop at NeurIPS*.
- [84] Mehrotra, A., Zampetakis, M., Kassianik, P., Nelson, B., Anderson, H., Singer, Y., & Karbasi, A. (2023). Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*.
- [85] Mei, A., Levy, S., & Wang, W. Y. (2023). Assert: Automated safety scenario red teaming for evaluating the robustness of large language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- [86] Mei, K., Fereidooni, S., & Caliskan, A. (2023). Bias against 93 stigmatized groups in masked language models and downstream sentiment classification tasks. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, (pp. 1699–1710).
- [87] Microsoft (2023). Planning red teaming for large language models (llms) and their applications - azure openai service.
URL <https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/red-teaming>
- [88] Mok, A. (2023). Chatgpt will no longer comply if you ask it to repeat a word 'forever'—after a recent prompt revealed training data and personal info.
URL <https://www.businessinsider.com/chatgpt-ai-refuse-to-respond-prompt-asking-repeat-word-forever-2023-12>
- [89] Mu, N., Chen, S., Wang, Z., Chen, S., Karamardian, D., Aljeraisy, L., Hendrycks, D., & Wagner, D. (2023). Can llms follow simple rules? *arXiv preprint arXiv:2311.04235*.
- [90] Nasr, M., Carlini, N., Hayase, J., Jagielski, M., Cooper, A. F., Ippolito, D., Choquette-Choo, C. A., Wallace, E., Tramèr, F., & Lee, K. (2023). Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*.

- [91] Neel, S., & Chang, P. (2023). Privacy issues in large language models: A survey. *arXiv preprint arXiv:2312.06717*.
- [92] Nguyen, C., Morgan, C., & Mittal, S. (2022). Poster cti4ai: Threat intelligence generation and sharing after red teaming ai models. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, (pp. 3431–3433).
- [93] Omrani Sabbaghi, S., Wolfe, R., & Caliskan, A. (2023). Evaluating biased attitude associations of language models in an intersectional context. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, (pp. 542–553).
- [94] OpenAI (2023). Gpt-4 technical report. (arXiv:2303.08774). ArXiv:2303.08774 [cs]. URL <http://arxiv.org/abs/2303.08774>
- [95] Peltine, K., Taufeeque, M., Zając, M., McLean, E., & Gleave, A. (2023). Exploiting novel gpt-4 apis. *arXiv preprint arXiv:2312.14302*.
- [96] Perez, E., Huang, S., Song, F., Cai, T., Ring, R., Aslanides, J., Glaese, A., McAleese, N., & Irving, G. (2022). Red teaming language models with language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, (pp. 3419–3448).
- [97] Pfau, J., Infanger, A., Sheshadri, A., Panda, A., Michael, J., & Huebner, C. (2023). Eliciting language model behaviors using reverse language models. In *Socially Responsible Language Modelling Research (SoLaR) Workshop at NeurIPS*.
- [98] Price, E. (2023). Asking chatgpt to repeat words “forever” may violate openai’s terms. URL <https://www.pcmag.com/news/asking-chatgpt-to-repeat-words-forever-may-violate-openais-terms>
- [99] Puttapparthi, P. C. R., Deo, S. S., Gul, H., Tang, Y., Shang, W., & Yu, Z. (2023). Comprehensive evaluation of chatgpt reliability through multilingual inquiries. *arXiv preprint arXiv:2312.10524*.
- [100] Qi, X., Huang, K., Panda, A., Wang, M., & Mittal, P. (2023). Visual adversarial examples jailbreak aligned large language models. In *The Second Workshop on New Frontiers in Adversarial Machine Learning*.
- [101] Qiu, H., Zhang, S., Li, A., He, H., & Lan, Z. (2023). Latent jailbreak: A benchmark for evaluating text safety and output robustness of large language models. *arXiv preprint arXiv:2307.08487*.
- [102] Radharapu, B., Robinson, K., Aroyo, L., & Lahoti, P. (2023). Aart: Ai-assisted red-teaming with diverse data generation for new llm-powered applications. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, (pp. 380–395).
- [103] Rae, J. W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., Aslanides, J., Henderson, S., Ring, R., Young, S., et al. (2021). Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.
- [104] Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., & Chen, M. (2022). Hierarchical text-conditional image generation with clip latents. (arXiv:2204.06125). ArXiv:2204.06125 [cs]. URL <http://arxiv.org/abs/2204.06125>

- [105] Rando, J., Paleka, D., Lindner, D., Heim, L., & Tramer, F. (2022). Red-teaming the stable diffusion safety filter. In *NeurIPS ML Safety Workshop*.
- [106] Rando, J., & Tramèr, F. (2023). Universal jailbreak backdoors from poisoned human feedback. *arXiv preprint arXiv:2311.14455*.
- [107] Rao, A., Vashistha, S., Naik, A., Aditya, S., & Choudhury, M. (2023). Tricking llms into disobedience: Understanding, analyzing, and preventing jailbreaks. *arXiv preprint arXiv:2305.14965*.
- [108] Rastogi, C., Tulio Ribeiro, M., King, N., Nori, H., & Amershi, S. (2023). Supporting human-ai collaboration in auditing llms with llms. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, (pp. 913–926).
- [109] Robey, A., Wong, E., Hassani, H., & Pappas, G. (2023). Smoothllm: Defending large language models against jailbreaking attacks. In *Robustness of Few-shot and Zero-shot Learning in Large Foundation Models (R0-FoMo) Workshop at NeurIPS*.
- [110] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, (p. 10674–10685). New Orleans, LA, USA: IEEE. URL <https://ieeexplore.ieee.org/document/9878449/>
- [111] Röttger, P., Kirk, H. R., Vidgen, B., Attanasio, G., Bianchi, F., & Hovy, D. (2023). Xstest: A test suite for identifying exaggerated safety behaviours in large language models. *arXiv preprint arXiv:2308.01263*.
- [112] Röttger, P., Vidgen, B., Nguyen, D., Waseem, Z., Margetts, H., & Pierrehumbert, J. (2021). Hatecheck: Functional tests for hate speech detection models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, (pp. 41–58).
- [113] Roy, S., Harshvardhan, A., Mukherjee, A., & Saha, P. (2023). Probing llms for hate speech detection: strengths and vulnerabilities. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, (pp. 6116–6128).
- [114] Roy, S. S., Naragam, K. V., & Nilizadeh, S. (2023). Generating phishing attacks using chatgpt. *arXiv preprint arXiv:2305.05133*.
- [115] Roy, S. S., Thota, P., Naragam, K. V., & Nilizadeh, S. (2023). From chatbots to phishbots?—preventing phishing scams created using chatgpt, google bard and claude. *arXiv preprint arXiv:2310.19181*.
- [116] Salem, A., Paverd, A., & Köpf, B. (2023). Maatphor: Automated variant analysis for prompt injection attacks. *arXiv preprint arXiv:2312.11513*.
- [117] Schlarmann, C., & Hein, M. (2023). On the adversarial robustness of multi-modal foundation models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, (pp. 3677–3685).
- [118] Schuett, J., Dreksler, N., Anderl jung, M., McCaffary, D., Heim, L., Bluemke, E., & Garfinkel, B. (2023). Towards best practices in agi safety and governance: A survey of expert opinion. *arXiv preprint arXiv:2305.07153*.

- [119] Schulhoff, S. V., Pinto, J., Khan, A., Bouchard, L.-F., Si, C., Anati, S., Tagliabue, V., Kost, A. L., Carnahan, C. R., & Boyd-Graber, J. L. (2023). Ignore this title and hackaprompt: Exposing systemic vulnerabilities of llms through a global prompt hacking competition. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- [120] Shah, R., Montixi, Q. F., Pour, S., Tagade, A., & Rando, J. (2023). Scalable and transferable black-box jailbreaks for language models via persona modulation. In *Socially Responsible Language Modelling Research (SoLaR) Workshop at NeurIPS*.
- [121] Shayegani, E., Dong, Y., & Abu-Ghazaleh, N. (2023). Jailbreak in pieces: Compositional adversarial attacks on multi-modal language models. (arXiv:2307.14539). ArXiv:2307.14539 [cs].
URL <http://arxiv.org/abs/2307.14539>
- [122] Shayegani, E., Mamun, M. A. A., Fu, Y., Zaree, P., Dong, Y., & Abu-Ghazaleh, N. (2023). Survey of vulnerabilities in large language models revealed by adversarial attacks. *arXiv preprint arXiv:2310.10844*.
- [123] Shen, X., Chen, Z., Backes, M., Shen, Y., & Zhang, Y. (2023). "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. *arXiv preprint arXiv:2308.03825*.
- [124] Shevlane, T., Farquhar, S., Garfinkel, B., Phuong, M., Whittlestone, J., Leung, J., Kokotajlo, D., Marchal, N., Anderljung, M., Kolt, N., et al. (2023). Model evaluation for extreme risks. *arXiv preprint arXiv:2305.15324*.
- [125] Shi, Z., Wang, Y., Yin, F., Chen, X., Chang, K.-W., & Hsieh, C.-J. (2023). Red teaming language model detectors with language models. *arXiv preprint arXiv:2305.19713*.
- [126] Solaiman, I. (2023). The gradient of generative ai release: Methods and considerations. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, (pp. 111–122).
- [127] Srivastava, A., Ahuja, R., & Mukku, R. (2023). No offense taken: Eliciting offensiveness from language models. *arXiv preprint arXiv:2310.00892*.
- [128] Stone, B., & Bergen, M. (2024). Openai working with u.s. military on cybersecurity tools.
URL <https://time.com/6556827/openai-us-military-cybersecurity/>
- [129] Sun, Z., Shen, Y., Zhou, Q., Zhang, H., Chen, Z., Cox, D., Yang, Y., & Gan, C. (2023). Principle-driven self-alignment of language models from scratch with minimal human supervision. *arXiv preprint arXiv:2305.03047*.
- [130] Tan, S., Joty, S., Baxter, K., Taeihagh, A., Bennett, G. A., & Kan, M.-Y. (2021). Reliability testing for natural language processing systems. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, (pp. 4153–4169).
- [131] The Frontier Model Forum (FMF) (2023). Frontier model forum: What is red-teaming?
URL <https://www.frontiermodelforum.org/uploads/2023/10/FMF-AI-Red-Teaming.pdf>
- [132] Tian, Y., Yang, X., Zhang, J., Dong, Y., & Su, H. (2023). Evil geniuses: Delving into the safety of llm-based agents. *arXiv preprint arXiv:2311.11855*.

- [133] Tong, S., Jones, E., & Steinhardt, J. (2023). Mass-producing failures of multimodal systems with language models. *arXiv preprint arXiv:2306.12105*.
- [134] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C. C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M.-A., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R., Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., & Scialom, T. (2023). Llama 2: Open foundation and fine-tuned chat models. (arXiv:2307.09288). ArXiv:2307.09288 [cs].
URL <http://arxiv.org/abs/2307.09288>
- [135] Tsai, Y.-L., Hsu, C.-Y., Xie, C., Lin, C.-H., Chen, J.-Y., Li, B., Chen, P.-Y., Yu, C.-M., & Huang, C.-Y. (2023). Ring-a-bell! how reliable are concept removal methods for diffusion models? *arXiv preprint arXiv:2310.10012*.
- [136] Tu, H., Cui, C., Wang, Z., Zhou, Y., Zhao, B., Han, J., Zhou, W., Yao, H., & Xie, C. (2023). How many unicorns are in this image? a safety evaluation benchmark for vision llms. *arXiv preprint arXiv:2311.16101*.
- [137] Wang, B., Chen, W., Pei, H., Xie, C., Kang, M., Zhang, C., Xu, C., Xiong, Z., Dutta, R., Schaeffer, R., et al. (2023). Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *arXiv preprint arXiv:2306.11698*.
- [138] Wang, H., & Shu, K. (2023). Backdoor activation attack: Attack large language models using activation steering for safety-alignment. *arXiv preprint arXiv:2311.09433*.
- [139] Wang, J., Wu, J., Chen, M., Vorobeychik, Y., & Xiao, C. (2023). On the exploitability of reinforcement learning with human feedback for large language models. *arXiv preprint arXiv:2311.09641*.
- [140] Wang, Y., Teng, Y., Huang, K., Lyu, C., Zhang, S., Zhang, W., Ma, X., & Wang, Y. (2023). Fake alignment: Are llms really aligned well? *arXiv preprint arXiv:2311.05915*.
- [141] Wang, Z., Yang, F., Wang, L., Zhao, P., Wang, H., Chen, L., Lin, Q., & Wong, K.-F. (2023). Self-guard: Empower the llm to safeguard itself. *arXiv preprint arXiv:2310.15851*.
- [142] Wei, A., Haghtalab, N., & Steinhardt, J. (2023). Jailbroken: How does llm safety training fail? In *Thirty-seventh Conference on Neural Information Processing Systems*.
- [143] Wei, Z., Wang, Y., & Wang, Y. (2023). Jailbreak and guard aligned language models with only few in-context demonstrations. *arXiv preprint arXiv:2310.06387*.
- [144] Weidinger, L., Rauh, M., Marchal, N., Manzini, A., Hendricks, L. A., Mateos-Garcia, J., Bergman, S., Kay, J., Griffin, C., Bariach, B., et al. (2023). Sociotechnical safety evaluation of generative ai systems. *arXiv preprint arXiv:2310.11986*.

- [145] Widder, D. G., West, S., & Whittaker, M. (2023). Open (for business): Big tech, concentrated power, and the political economy of open ai. (4543807).
URL <https://papers.ssrn.com/abstract=4543807>
- [146] Wood, B. J., & Duggan, R. A. (2000). Red teaming of advanced information assurance concepts. In *Proceedings DARPA Information Survivability Conference and Exposition. DIS-CEX'00*, vol. 2, (pp. 112–118). IEEE.
- [147] Wu, F., Liu, X., & Xiao, C. (2023). Deceptprompt: Exploiting llm-driven code generation via adversarial natural language instructions. *arXiv preprint arXiv:2312.04730*.
- [148] Wu, Y., Li, X., Liu, Y., Zhou, P., & Sun, L. (2023). Jailbreaking gpt-4v via self-adversarial attacks with system prompts. *arXiv preprint arXiv:2311.09127*.
- [149] Xie, Y., Yi, J., Shao, J., Curl, J., Lyu, L., Chen, Q., Xie, X., & Wu, F. (2023). Defending chatgpt against jailbreak attack via self-reminders. *Nature Machine Intelligence*, (pp. 1–11).
- [150] Xu, J., Ju, D., Li, M., Boureau, Y.-L., Weston, J., & Dinan, E. (2020). Recipes for safety in open-domain chatbots. *arXiv preprint arXiv:2010.07079*.
- [151] Xu, J., Ju, D., Li, M., Boureau, Y.-L., Weston, J., & Dinan, E. (2021). Bot-adversarial dialogue for safe conversational agents. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, (pp. 2950–2968).
- [152] Xu, N., Wang, F., Zhou, B., Li, B. Z., Xiao, C., & Chen, M. (2023). Cognitive overload: Jailbreaking large language models with overloaded logical thinking. *arXiv preprint arXiv:2311.09827*.
- [153] Yang, Y., Hui, B., Yuan, H., Gong, N., & Cao, Y. (2023). Sneakyprompt: Jailbreaking text-to-image generative models. (arXiv:2305.12082). ArXiv:2305.12082 [cs].
URL <http://arxiv.org/abs/2305.12082>
- [154] Yao, D., Zhang, J., Harris, I. G., & Carlsson, M. (2023). Fuzzllm: A novel and universal fuzzing framework for proactively discovering jailbreak vulnerabilities in large language models. *arXiv preprint arXiv:2309.05274*.
- [155] Yao, Y., Duan, J., Xu, K., Cai, Y., Sun, E., & Zhang, Y. (2023). A survey on large language model (llm) security and privacy: The good, the bad, and the ugly. *arXiv preprint arXiv:2312.02003*.
- [156] Yong, Z. X., Menghini, C., & Bach, S. (2023). Low-resource languages jailbreak gpt-4. In *Socially Responsible Language Modelling Research (SoLaR) Workshop at NeurIPS*.
- [157] Yu, J., Lin, X., & Xing, X. (2023). Gptfuzzer: Red teaming large language models with auto-generated jailbreak prompts. *arXiv preprint arXiv:2309.10253*.
- [158] Yuan, Y., Jiao, W., Wang, W., Huang, J.-t., He, P., Shi, S., & Tu, Z. (2023). Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher. *arXiv preprint arXiv:2308.06463*.
- [159] Zenko, M. (2015). *Red Team: How to succeed by thinking like the enemy*. Basic Books.

- [160] Zhang, J., Zhou, Y., Hui, B., Liu, Y., Li, Z., & Hu, S. (2023). Trojansql: Sql injection against natural language interface to database. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, (pp. 4344–4359).
- [161] Zhang, M., Pan, X., & Yang, M. (2023). Jade: A linguistics-based safety evaluation platform for llm. *arXiv preprint arXiv:2311.00286*.
- [162] Zhang, X., Zhang, C., Li, T., Huang, Y., Jia, X., Xie, X., Liu, Y., & Shen, C. (2023). A mutation-based method for multi-modal jailbreaking attack detection. *arXiv preprint arXiv:2312.10766*.
- [163] Zhang, Z., Yang, J., Ke, P., & Huang, M. (2023). Defending large language models against jailbreaking attacks through goal prioritization. *arXiv preprint arXiv:2311.09096*.
- [164] Zhao, W., Li, Z., & Sun, J. (2023). Causality analysis for evaluating the security of large language models. *arXiv preprint arXiv:2312.07876*.
- [165] Zhu, S., Zhang, R., An, B., Wu, G., Barrow, J., Wang, Z., Huang, F., Nenkova, A., & Sun, T. (2023). Autodan: Automatic and interpretable adversarial attacks on large language models. In *Socially Responsible Language Modelling Research (SoLaR) Workshop at NeurIPS*.
- [166] Zhu, Z., Wang, J., Cheng, H., & Liu, Y. (2023). Unmasking and improving data credibility: A study with datasets for training harmless language models. *arXiv preprint arXiv:2311.11202*.
- [167] Zhuo, T. Y., Huang, Y., Chen, C., & Xing, Z. (2023). Red teaming chatgpt via jailbreaking: Bias, robustness, reliability and toxicity. *arXiv preprint arXiv:2301.12867*, (pp. 12–2).
- [168] Zou, A., Wang, Z., Kolter, J. Z., & Fredrikson, M. (2023). Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

A Research Survey and Case Study Analysis Details

This appendix contains further details about the research papers and case studies explored as part of this work. Table 3 contains the specific classifications along each dimension described in Section 4 for every work recovered as part of our research survey. We additionally provide access to a Google Sheets project with notes and thematic analyses for both the case studies and research papers retrieved and described in this work. The project can be accessed via this URL: <https://docs.google.com/spreadsheets/d/1cZPc6A1kf8sq0FMsEvZgI2PzX2tHTIbemMa6sq4J2Qk/edit?usp=sharing>.

		Type of Risk Investigated			
Type of Approach Used		Subjective	Objective	Both	Neither
		Brute-Force	[167, 151, 150, 112, 53, 45, 119, 60, 76, 113, 143, 123, 107, 79, 149]	[36, 89, 114, 75]	[51]
Brute-Force + AI	[101, 129, 71, 82, 16, 157, 35, 40, 85, 15, 83, 120, 97, 111, 163, 43, 140, 141, 34, 154, 63, 9, 161]	[25, 125, 56, 147, 166, 136, 115]	[96, 108, 127, 142, 57, 148, 74, 116, 132, 137]	[28, 102]	
Algorithmic Search	[37, 81, 27, 135, 165, 52, 33, 84, 109, 23, 69, 153]	[92]	[162]	None	
Targeted Attack	[105, 156, 61, 164, 24, 139, 138, 158, 168, 70, 99, 78, 152, 55, 42, 121, 41, 100, 106]	[160, 133, 13, 117, 26, 90, 6]	[95, 20]	None	

Table 3: In-depth classification of papers acquired for our survey based on the type of content produced and type of approach used in each paper. See Section 4 for details and definitions.

This figure "acm-jdslogo.png" is available in "png" format from:

<http://arxiv.org/ps/2401.15897v1>