# 🧠 Operationalizing a Threat Model for Red-Teaming Large Language Models (LLMs)

**Anonymous authors**
**Paper under double-blind review**

## Abstract

Creating secure and resilient applications with large language models (LLM) requires anticipating, adjusting to, and countering unforeseen threats. Red-teaming has emerged as a critical technique for identifying vulnerabilities in real-world LLM implementations. This paper presents a detailed threat model and provides a systematization of knowledge (SoK) of red-teaming attacks on LLMs. We develop a taxonomy of attacks based on the stages of the LLM development and deployment process and extract various insights from previous research. In addition, we compile methods for defense and practical red-teaming strategies for practitioners. By delineating prominent attack motifs and shedding light on various entry points, this paper provides a framework for improving the security and robustness of LLM-based systems.

## 1 Introduction

> *Red Teaming is the process of using tactics,*
> *techniques, and processes (TTP) to emulate a*
> *real-world threat with the goal of training and*
> *measuring the effectiveness of people, processes,*
> *and technology used to defend an environment.*
>
> Vest & Tubberville (2020)

The practice of "red-teaming" originated during the 1960s in United States (US) military Cold War simulations to anticipate threats from Soviet Union (Averch & Lavin, 1964; Red Team). In a red-teaming exercise, the "red team" plays the role of an adversary and attempts to compromise a system, while a blue team plays the role of a defender and is responsible for fixing the security gaps in the system. The purpose of red-teaming is to adopt an adversarial mindset to identify weaknesses and security vulnerabilities in a system. Over time, red-teaming expanded beyond military operations to domains such as cybersecurity (Duggan & Wood, 1999), airport security (Price, 2004), and more recently to AI and machine learning (ML), and specifically to large language models (LLMs) (OpenAI, 2023) and generative AI.

Paradoxically, LLMs are predictable and unpredictable. On the one hand, these models are highly predictable because the lower model loss of a better LLM means that it performs better in predicting the next words. On the other hand, they are also highly unpredictable, as the universality of being adept at predicting the next words makes it impossible to anticipate the specific capabilities and outputs of the LLM a priori before fine-tuning it on the downstream task Wei et al. (2022); Ganguli et al. (2022a). This unpredictability poses a challenge in understanding the consequences of deploying LLMs in real-world scenarios. For example, LLMs can hallucinate (Verma & Oremus, 2023; Dale et al., 2023), reveal personally identifiable information (PII) (Carlini et al., 2020a), be used to generate misinformation (Hazell, 2023; Krishna et al., 2024), generate biased (Santurkar et al., 2023; Hartmann et al., 2023; Ghosh & Caliskan, 2023; Sabbaghi et al., 2023), unsettling (Kevin, 2023a;b), sycophantic (Perez et al., 2023), toxic (Perez et al., 2022; Shen et al., 2023c), harmful (Weidinger et al., 2021) and insecure (Majdinasab et al., 2023) content in response to benign or malicious

prompts. A lack of transparency surrounding the development of these models further exacerbates these issues (Widder et al., 2023; Zhang et al., 2024b; Casper et al., 2024b).

Red-teaming has emerged as an essential tool for assessing the safety of LLMs and minimizing the risks associated with their deployment in human-facing products. To this end, industry and academia have developed and published approaches in red teaming in ML (Pearce & Lucas, 2023; Fabian, 2023; Siva Kumar, 2023b; Meta, 2023; Ji, 2023). Competitions such as DEFCON (Cattell et al., 2023) and the RLHF Trojan Competition (Rando & Tramèr, 2024a;b), along with games such as Hacc-man (Valentim et al., 2024) and Tensor Trust (Toyer et al., 2024), as well as recently published regulations (Bletchley Declaration, 2023) and U.S. Executive Orders (The White House, 2023) have led to widespread awareness of this topic and shed light on red-teaming strategies.

Through this paper, our objective is to systematize knowledge about the current state of LLM red-teaming, allowing researchers and practitioners to navigate the complexities of developing **H**elpful, **H**armless, and **H**onest LLM-based applications (**H³LLM**) (Askell et al., 2021). Drawing on previous research efforts, we develop a taxonomy of red-teaming attacks to summarize various aspects of this emerging field (see Figure 1). Our main contributions are as follows.

1. We introduce a threat model based on entry points in the LLM development and deployment lifecycle to allow reasoning about various kinds of attack and associated defenses.

2. We provide a taxonomy of attacks based on our proposed threat model, followed by a brief discussion of common defense methodologies.

3. Finally, we systematize various insights derived from previous published works to tease out the desirable properties needed to conduct effective red-teaming exercises and ensure robust defense strategies.

**Paper Outline:** In Section § 2, we provide an overview of LLM training and inference phases and define harmful behaviors. We then draw a distinction between red-teaming and traditional evaluations of trustworthiness and bias. In Section § 3 we describe the threat model and the fundamental principle that we use to structure various attacks in a taxonomy. In Section § 4, we summarize various attack methods followed by briefly addressing defenses in Section § 5. We provide a discussion of the insights gathered from previous work in Section § 6. And finally, we conclude our study by distilling the current state of red-teaming to identify several promising future research avenues in Section § 7.

## 2 Background: LLMs, Harms, and the Red-Teaming Paradigm

In this section, we provide details on the LLM life cycle from training to deployment. We define harmful behavior and draw a distinction between red-teaming and traditional fairness evaluations. Finally, we highlight the advantage of systematizing attacks based on the proposed threat model compared to previous surveys.

### 2.1 LLM Development Phases

**Pre-training:** In this initial stage, LLMs learn from a large dataset, acquiring basic language understanding and context Biderman et al. (2023); Brown et al. (2020b).

**Supervised Fine-Tuning (SFT) / Instruction Tuning (IT):** After pre-training, models are fine-tuned with specific datasets to adapt to particular tasks or domains. Fundamentally, SFT helps LLMs understand the semantic meaning of prompts and produce relevant responses Ouyang et al. (2022b); Lou et al. (2023).

**Reinforcement Learning from Human Feedback (RLHF):** This phase involves the refinement of the model responses based on human feedback, focusing on aligning the outputs with human values and expectations (Rafailov et al., 2023; Ouyang et al., 2022a; Rafailov et al., 2023; Bai et al., 2022; Korbak et al., 2023). For brevity, we omit an extended discussion of alignment but refer to Shen et al. (2023b) and Wang et al. (2023e) for a comprehensive overview.

**Deployment:** Finally, a trained LLM can be deeply integrated in consumer technology applications like chatbots, email, code review, news summarization, and legal document analysis, among other uses with unmediated access to surrounding components (Yang et al., 2023a).

## 2.2 Defining Harmful Behavior

Defining harmful behavior for a red-teaming exercise requires a nuanced understanding of "harm." Acceptable use policies of OpenAI, Anthropic, Inflection AI, Perplexity AI, among others, can serve as a good starting point for defining harmful behavior (Usage Policy OpenAI, 2023; Usage Policy Anthropic, 2023). Following the NIST Common Vulnerabilities and Exposures (CVE) guidelines (NIST CVE, 2022), the Avid Taxonomy Matrix (Avid Taxonomy Matrix, 2023) provides a holistic taxonomy of risk categories and failure modes associated with an ML system. In addition to conventional security vulnerabilities, the Open Web Application Security Project (OWASP) has also published top 10 risks for GenAI applications (OWASP, 2025).

| Study | Description |
| --- | --- |
| (Inan et al., 2023) | Llama Guard Taxonomy |
| (Wang et al., 2023a) | Decoding Trust Taxonomy |
| (OpenAI, 2024b) | OpenAI Moderation Endpoint Categories |
| (Vidgen et al., 2024) | AI Safety Benchmark Hazard Categories |
| (Tedeschi et al., 2024) | ALERT Risk Taxonomy |
| (Perspective API, 2024) | Perspective API Toxicity Attributes |
| (Avid Taxonomy Matrix, 2023) | Avid Taxonomy Matrix |
| (Ghosh et al., 2024) | AEGIS Risk Taxonomy |
| (Zeng et al., 2024a) | AI Risk Categorization |

Table 1: Sample risk taxonomies that can guide practitioners in developing domain-specific or application-specific risk taxonomies

Table 1 illustrates examples of several risk taxonomies. Some common risk categories in these taxonomies are Sexual Content, Violence & Hate, etc. Feffer et al. (2024) categorize these risk categories into two broad buckets - dissentive risks and consentive risks. Dissentive risks comprise risk categories whose definitions are not widely agreed upon (for example, some people might find the response to "how to build a bomb?" admissible), while consentive risks are risks whose definition is widely agreed upon and no additional context is needed (for example, data and private information leakage, phishing attacks are inadmissible in any context) (Feffer et al., 2024). Depending on the particular domain, practitioners may need to expand these risk taxonomies to include domain-specific risk categories. (e.g., investment advice abstention for financial domain, citation on point for legal domain, etc. (Tshimula et al., 2024)) Furthermore, certain behaviors could be of *dual intent* (Mazeika et al., 2024; Stapleton et al., 2023). For example, writing encryption functions could be performed by cyber-security developers or by malicious hackers. Similarly, in some cases, it may be important to focus on the *differential harm* caused by an LLM over online searchability to quantify the additional harm introduced for the given input (Mazeika et al., 2024). It is worth noting that harmful behavior can arise without malicious intent (e.g., hallucinations).
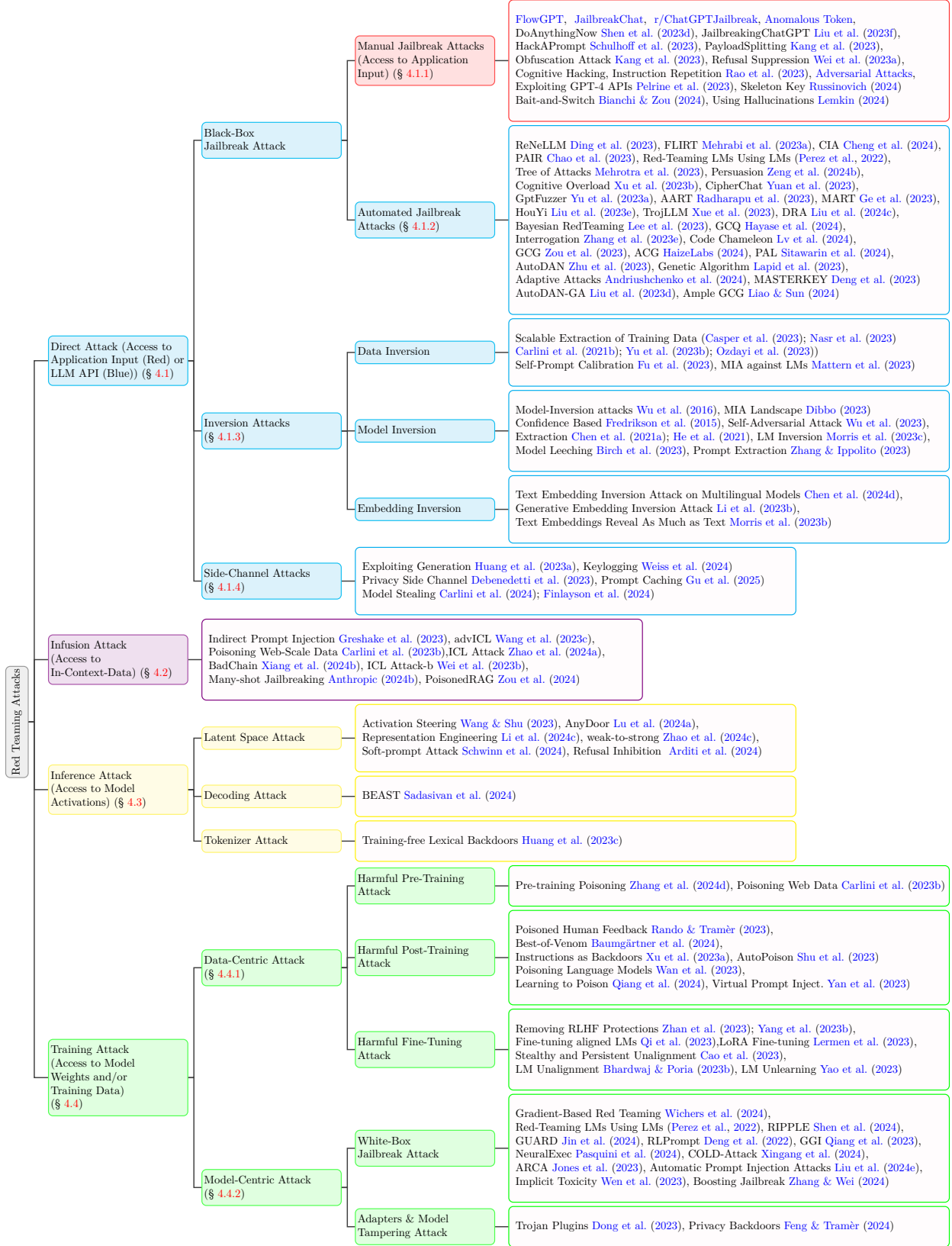
Figure 1: Taxonomy of LLM red-teaming attacks, ranging from prompt-based to training-level attacks based on required access levels.

Additionally, emerging threats and harms, such as the use of LLMs for targeted influence operations (Goldstein et al., 2023), require more attention and input from subject matter experts to properly define the harm. Previous work emphasizes the need to clearly define the risks and behaviors to uncover before starting red-teaming (Casper et al., 2023; Feffer et al., 2024).

**Context-Dependent Nature of Harm.**   The concept of harm in LLM outputs is highly context-dependent. Creative hallucinations, for example, might be beneficial for a writer seeking inspiration, but could be detrimental in scenarios requiring factual accuracy, such as legal advice or medical information Tamkin et al. (2021). The working definition of harm will vary depending on the specific domain (e.g., laws, rules, and social norms) in which it is used.

**Scope of Harmful Behavior.**   This paper narrows its focus to the types of harmful behavior that are particularly relevant for LLM applications. Ferrara (2023); Greshake et al. (2023) provide an overview of the nefarious uses of LLMs. For a broader discussion that encompasses bias, fairness, legal and regulatory considerations, the reader is directed to the foundational literature in the field that describes specific definitions and frameworks for understanding and reducing harm in ML systems Ding et al. (2021); Casper et al. (2024a).

## 2.3   Safety vs Security Objectives

Safety and security represent distinct objectives in LLM risk management (Qi et al., 2024a). Safety focuses on preventing harm that LLMs might inflict upon their environment - for example, generating toxic content or providing harmful advice accidentally. Security focuses on protecting LLMs against malicious exploitation - such as preventing jailbreaking attacks or unauthorized access to training data.

These objectives entail different threat models. Safety primarily addresses non-adversarial scenarios like unintended model behaviors and inherent flaws. Security explicitly considers adversarial scenarios where malicious actors attempt to compromise the system. For example, hallucination is primarily a safety concern when occurring naturally but becomes a security issue when adversaries deliberately induce it through attacks.

In this paper, we adopt a security-focused perspective by considering an adversarial threat model where bad actors actively attempt to compromise LLM systems while ML practitioners aim to defend against such attacks. Our threat model systematically analyzes attack surfaces and entry points that adversaries may exploit, from jailbreaking attempts at the application layer to more sophisticated attacks targeting model weights and training processes. This security-oriented approach complements existing safety research while specifically focusing on defending against malicious exploitation.

## 2.4   Benchmark Evaluation and Red-Teaming Evaluation

Conventional evaluation methods are generally divided into two main categories: (1) Automatic Evaluation and (2) Human Evaluation. Automatic evaluation involves comparing model outputs with ground-truth annotations obtained offline to calculate a metric. It may also include having a more advanced model assess the accuracy of the model's output when ground truth references are unavailable or when the task is too open-ended to be accurately measured by a limited set of references (e.g., text summarization). Human evaluation, on the other hand, uses human judges to assess the accuracy or quality of a model's output. This type of evaluation can also be designed for contrastive assessment, where a human rater chooses their preferred output between two model predictions (e.g., learning reward models in RLHF methods). Evaluation helps answer questions about model capabilities, the effectiveness of the training algorithm, and provides a way to compare different models. Various benchmarks have been proposed to assess the trustworthiness and safety of an LLM (Wang et al., 2023a; Huang et al., 2023b; Sun et al., 2024).

Unlike conventional evaluation and benchmarking practices focused on measuring performance and fairness, red-teaming adopts a proactive approach aimed at uncovering potential vulnerabilities that can lead to catastrophic failure. Inie et al. (2023) define red-teaming as, "A limit-seeking activity, using vanilla attacks, a manual process, team effort, and an alchemist mindset to break, probe, or experiment with LLMs." while Barrett et al. (2024) describes it as, "more intensive and interactive testing by domain experts, providing a deeper understanding of a model's behavior in various scenarios." Using a Software quality assurance (SQA)

analogy, evaluation is like running unit and regression tests, while red-teaming is like discovering bugs and writing new test cases for them.

By simulating adversarial attacks, practitioners can better understand and fortify models against misuse in the real world. Robust red-teaming can help facilitate ethical and safe applications in various contexts Nwadike et al. (2020). However, as Feffer et al. (2024) state "red-teaming is not a panacea" and should complement other forms of ML governance and model evaluation (Shevlane et al., 2023).

## 2.5 Related Work

**User Surveys and Interviews:** Inie et al. (2023) describe the insights from user interviews to understand the red-teaming mindset and the primary motivations behind participating in such an activity. They state, "The primary motivations for partaking in the activity were curiosity, fun, and concerns (intrinsic), and professional and social (extrinsic)" (Inie et al., 2023). Similarly, Schuett et al. (2023) conducted a survey to collect expert opinion on the role that leading AI laboratories should play in AI safety. They discovered that 98% of the participants concurred that AGI laboratories should evaluate risks prior to deployment, examine hazardous capabilities, perform third-party audits, enforce usage limitations, and engage in red-teaming.

**LLM Attack Taxonomies:** In previous red-teaming surveys, attacks are primarily organized using two modes. The first group organizes attacks based on the type of risk posed, such as hate speech, misinformation, and privacy leakage. These categories emerge naturally from the risk taxonomy described above (see Section § 2.2). Representative works in this group include Weidinger et al. (2021); Abdali et al. (2024); Greshake et al. (2023); Zhuo et al. (2023); Wang et al. (2023a). This way of organizing is more useful to a policy maker trying to assess the readiness of the model across risk categories than to a practitioner trying to identify specific failure points in model development lifecycle.

The second group of work organizes attacks based on the methodology used for the attack (e.g., automatic, manual, etc.). Representative works here are Inie et al. (2023); Chowdhury et al. (2024); Lin et al. (2024); Schulhoff et al. (2023); Dong et al. (2024b); Geiping et al. (2024); Feffer et al. (2024); ATLAS Matrix (2023); Shayegani et al. (2023). Of these, Schulhoff et al. (2023) and Geiping et al. (2024) are limited to prompt-based attacks and would correspond to the Black-Box Jailbreak Attack category in our taxonomy presented later (see Section § 4.1.1, § 4.1.2), while Dong et al. (2024b) broadly categorizes attacks into Training-Time or Inference-Time Attacks and Chowdhury et al. (2024) organizes attacks by three attack techniques, namely Jailbreaks, Prompt Injection, and Data Poisoning. ATLAS Matrix (2023) outlines 14 attack strategies mainly from the point of view of a security researcher, while Shayegani et al. (2023) classifies attacks according to the input modality. Although helpful, these organization schemes miss the nuances of various types of attacks and lack the description of a corresponding threat model.

To our knowledge Feffer et al. (2024) from the second group is most similar to our work, which categorizes red-teaming attacks into four broad categories: (1) Brute-force (2) Brute-force + AI (3) Algorithmic Search and (4) Targeted Attack. The last category "Targeted Attack" involves "deliberately targeting part of an LLM" (Feffer et al., 2024) and would be similar to what our proposed attack taxonomy aims to achieve.

Going beyond existing surveys, we offer an overview of attacks spanning multiple stages of the model life cycle, ranging from training to deployment, and organize attacks based on the level of access required to execute them. This method aligns more closely with the terminology used by machine learning professionals. Finally, we explore methods to counter these attacks and propose strategies for efficient red-teaming and defense. By organizing attacks around entry points and also linking them back to the threat model, we see a more complete picture of an adversary's capabilities and goals than has been seen in the specific prior literature.

**Scope:** This study does not cover vulnerabilities related to programming languages and cybersecurity exploits (Siva Kumar, 2023a; Sanseviero, 2024; CISA, 2024; Zhang et al., 2024a). Additionally, covert malware in AI development platforms (Goodin, 2024), watermark evasion attacks (Liu et al., 2023a) and attacks targeting vision-language models are outside the scope of this work (Liang et al., 2024; Yang et al., 2023c; Niu et al., 2024; Gong et al., 2023b).

# 3 Threat Model

A threat model refers to the potential vulnerabilities or risks for which a model is evaluated and the actions and information that the adversary has at their disposal to conduct an attack. By understanding how users interact with LLMs, how LLMs are trained, tested, and deployed, a clearer picture of the range of possible attacks emerges. In this section, we explore the nuances of user interaction through prompting, the implications of application layers beyond the core model, and the importance of defining and mitigating harmful behavior within diverse contexts.

**(1) Interaction through Prompting.** Human interaction with LLMs occurs primarily through prompting (i.e.,"question-answering"). Prompts mimic how humans interact in the non-digital world, permitting a wide range of flexible applications.

Some successful applications include chatbots [ChatGPT (OpenAI, 2022) and OpenAssistant (Kopf et al., 2023)], Voice Assistants (Soltan et al., 2022), code generation (Chen et al., 2021b), interactive fiction (Calderwood et al., 2022), news (Zhang et al., 2023a), medical (Tang et al., 2023) and judgement summarization (Deroy et al., 2023). However, prompts provide the flexibility that opens the door to various attacks aimed at eliciting harmful or unintended responses from an LLM. Understanding and mitigating these risks is essential for safe LLM deployment Bommasani et al. (2021).

**(2) Application Layers beyond LLMs.** The complexity of LLM applications often extends beyond the model itself, incorporating additional components such as retrieval systems, heuristic filters, and error correction mechanisms. LLMs can also be used in multistep planning and reasoning agents (Yao et al., 2022; Wang et al., 2023d; Zaharia et al., 2024; LLM Agents, 2023; Hamilton, 2023; Ziems et al., 2023; Wu et al., 2024c; Dasgupta et al., 2023), to call external tools (Schick et al., 2023; Patil et al., 2023), and collaborate with other agents to complete complex tasks. Due to this complexity, effective red-teaming must encompass the entire end-to-end application to fully assess vulnerabilities and protect against possible misuse (Fang et al., 2024b;a).

**(3) Model Internals and Training Data.** As described in Section §2, training LLMs involves several steps, from collecting web-scale datasets to preference and instruction tuning. Each of these steps opens a potential entry point for attacks.

## 3.1 Attack Surface

Understanding the attack surface of a system informs where and how an adversary may attempt to subvert that system. Figure 2 illustrates various entry points corresponding to different attacks in our taxonomy (Figure 1).

Attack methods are grouped according to the level of system access required. Jailbreak attacks require narrow access to just the LLM based application input. Training time attacks require a much wider access including access to training procedure and some subset of instruction tuning or fine-tuning data. In an offline training time attack scenario, an adversary poisons a subset of training data to insert backdoor attack phrases or alter the training algorithm to induce harmful behavior from the model. In the online training time attack scenario, an adversary manipulates model activations or corrupts the tokenizer. Boxes ④ and ⑤ encapsulate these attack vectors. Since these attacks require white-box access to model artifacts and training data, they are more likely to be executed by an insider (e.g., red team ML researcher) or a state-level actor (e.g., intelligence agencies). The human in the loop annotation process, not shown in the figure, is another potential point of vulnerability for poisoning instruction tuning or preference pair datasets.

Box ③ encapsulates attack vectors that arise from models using parametric (stored in model weights) and non-parametric (stored in external systems, e.g., retrieval systems, search indices, web pages, in-context examples, and databases) information to generate their output.

Boxes ② and ① represent black-box access to a model. Box ② encapsulates the LLM API parameters, proxy models, or vector embeddings which could be used to conduct a Direct Attack. The *logit_bias*
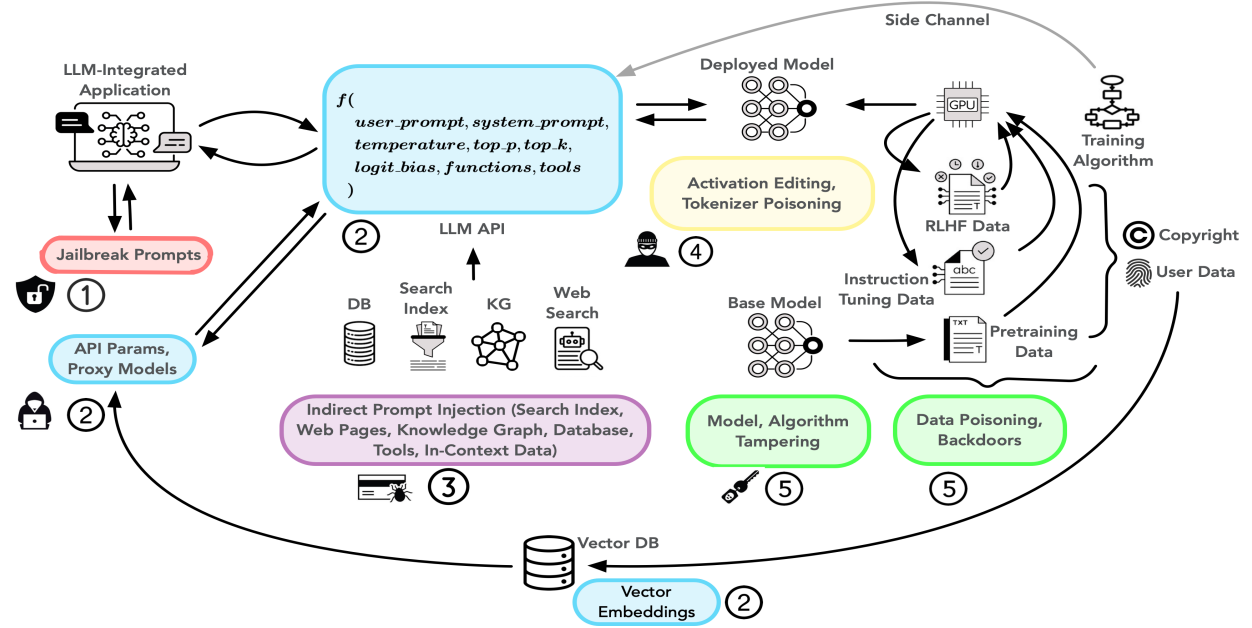
Figure 2: Attack vectors corresponding to the various attacks in our proposed taxonomy. Attacks are arranged in increasing order with respect to the level of access required. Attacks on the left side target late entry points in the lifecycle stage, such as application input, whereas attacks on the right side target early entry points such as training data, algorithm, etc. The colored boxes indicate the attack vectors corresponding to each high-level attack type in our taxonomy. Black arrows indicate the flow of information or artifacts, while gray arrows indicate side channels exposed due to the knowledge of common data-preprocessing steps such as data deduplication, etc. Boxes ② and ① represent black-box access to the model. Box ② encapsulates the LLM API params, vector embeddings and proxy models as the attack vectors. Box ③ encapsulates additional retrieval systems, tools, and in-context data as attack vectors. The boxes ④ and ⑤ represent white-box access to the model and encapsulate full or partial access to the model weights, the training algorithm, and the training data

parameter, which can be exploited to reveal token-level probabilities, is particularly notable. Box ① represents the user input in an LLM-integrated application. This corresponds to the Manual Jailbreak Attack in our taxonomy and represents the widest attack surface and can be easily executed by anyone. For instance, an LLM-provider interested in defending against bad press coverage from a hostile journalist is probably most concerned about Manual Jailbreak attacks.

In addition, system-level components, such as training data deduplication and output filtering, expose side channels that can leak information about the training data.

## 3.2 Adversary Capabilities

An adversary can be internal (member of the red-team, participant in the red-teaming competition), a weak eavesdropper (hobbyist, hacker, user), or a strong state-level adversary, among others. An adversary can inject, modify, or delete training data, or inject prompts into non-parametric knowledge sources, such as databases and search indices. They could also tamper with the learning algorithm or manipulate the activations of the model. Finally, an adversary can exploit the generation parameters, side channels, proxy models, embeddings, or prompts to leak sensitive information, discover harmful prompts, and induce harmful outputs.

An adversary's capabilities vary based on their level of access to the system components shown in Figure 2. Table 3.2 details the specific capabilities available to adversaries in each attack vector. Moving from Box 1 to Box 5, we observe an increasing sophistication in adversary capabilities, from simple prompt crafting to full

| Attack Vector | Attack Type | Capabilities |
|---|---|---|
| Application Input (Box 1) | Manual Attack | Craft malicious prompts, exploit application context, manipulate system instructions (Non-programmatic access) |
| LLM API (Box 2) | Direct Attack | Access generation parameters (temperature, top-k, logit-bias, etc.) (Programmatic access), exploit side channels (e.g., data processing or de-duplication filters), extract model probabilities through API parameters |
| In-Context Data (Box 3) | Infusion Attack | Poison retrieved documents, search-index, web-pages, manipulate in-context examples, compromise external tools and APIs which are invoked by an LLM |
| Model Internals (Box 4) | Inference Attack | Modify model activations, access embedding space, control decoding strategy, tokenizer tampering |
| Training Process (Box 5) | Training Attack | Poison training data, modify training algorithm, insert backdoors during model training, access to sufficient computational resources, erode model alignment through fine-tuning, learn adversarial prompts through compute-intensive prompt-search algorithms |

Table 2: Adversary Capabilities by Attack Vector

control over the training process. This progression of capabilities also typically correlates with increasing technical expertise and resources required to execute attacks effectively.

### 3.3 Adversary Goals

Finally, to complete our description of the threat model, we specify the goals of an adversary. For example, a hostile journalist might try to induce hallucinations or obtain factually inaccurate statements from a language model, while a malicious state-level actor could try to poison datasets at web scale or extract the parameters of black-box models to create rival services. A hacker could also attempt to insert phishing messages into the responses generated by these models or undermine code-generating language models by introducing software vulnerabilities in their output.

These goals can be modeled through the lens of Confidentiality, Integrity, Availability, and Privacy (CIAP) (Papernot et al., 2018). Referring to the analysis in Papernot et al. (2018) and extending it to LLMs, an adversary may attempt to extract model weights, application prompts, and other intellectual property (Targeting Confidentiality); exfiltrate sensitive user data or other confidential data that was used in model training (Targeting Privacy); attempt to generate harmful or incorrect outputs (Targeting Integrity); finally, attempt to degrade output quality, generation time or access (e.g., denial of service (Model Denial of Service, 2023), exhausting GPU resources, hitting API quotas) (Targeting Availability). Table 3.3 outlines the strategic goals that adversaries aim to achieve through different attack vectors, highlighting how system access enables increasingly sophisticated attacks on model behavior and security.

## 4 Attacks

With the threat model in mind, we operationalize the attack strategies and methods in order of access required from the least to the greatest. While there are many ways to categorize attacks, such as based on the type of approach used in constructing the attack or the type of risk being targeted by the attack, we organize attacks based on attack entry points, as it offers a clear understanding of adversary capabilities and allows practitioners to focus defense efforts on the most vulnerable points. In this section, we will cover various types of attack and discuss defenses in the next section. (see Section § 5). The sublevels within our taxonomy are grouped more loosely. For instance, Data Inversion Attacks under Inversion Attacks (Section

| Attack Vector | Strategic Goals |
|---|---|
| Application Input | Generate harmful content, bypass safety alignment, obtain unauthorized information through social engineering, damage reputation |
| LLM API | Extract training data or model-weights through API parameters, infer model capabilities, exploit side channels, build replica service |
| In-Context Data | Compromise downstream applications through poisoned retrieval, manipulate model behavior through contaminated examples, spread misinformation at scale |
| Model Internals | Control model outputs through activation engineering, extract information from embedding space |
| Training Process | Insert effective, stealthy and persistent backdoors for long-term targeted control, compromise alignment through adversarial training, corrupt model behavior, generate and share adversarial prompts |

Table 3: Strategic Goals Behind Different Attack Vectors

§ 4.1.3) are not the only attacks that can leak training data. It is also worth drawing a distinction between the planning and execution phases of an attack. For example, some backdoor attacks involve strategically implanting backdoor tokens in the training data during the planning phase and using these backdoor triggers to prompt a trained model during the execution phase to materialize the attack. Our taxonomy organizes attacks based on the highest level of access required during either the planning or execution phases of an attack. Analogously, in the White-Box Prompt Search Attack (refer to Section § 4.4.2), the planning phase requires executing a costly adversarial prompt search algorithm, whereas the execution phase entails directly prompting the LLM with the identified adversarial prompt.

## 4.1 Direct Attack

Direct Attacks encompass vulnerabilities that can be exploited through interaction with the model's external interfaces. We categorize these into three distinct classes: (1) Black-Box Jailbreak Attacks, which attempt to circumvent model safeguards, (2) Inversion Attacks, which aim to extract protected information, and (3) Side-Channel Attacks, which exploit architectural vulnerabilities. Through the lens of the CIAP framework (Papernot et al., 2018), these attacks target different security objectives: Jailbreak Attacks primarily compromise system Integrity, while Inversion and Side-Channel Attacks threaten Confidentiality, Privacy, and Availability. Within Black-Box Jailbreak Attacks, we make a crucial distinction between Manual and Automated approaches. Manual Attacks require only access to the application interface and can be executed without technical expertise. In contrast, Automated Attacks leverage programmatic API access, providing additional degrees of freedom through exposed API parameters (e.g., temperature, top-k sampling), stored embeddings, and proxy models (see box ②️ in Figure 2). This expanded attack surface, combined with the ability to systematically probe model behavior, enables more sophisticated attack strategies.

Notably, many sophisticated Automated Attacks trace their lineage to simpler Manual Attacks discovered by the broader AI safety community. For instance, the "Do Anything Now" (DAN) jailbreak (Shen et al., 2023d), originally shared on Reddit by an AI safety enthusiast, inspired more systematic approaches like AutoDAN (Zhu et al., 2023) and AutoDAN-GA (Liu et al., 2023d). Similarly, manual refusal suppression techniques pioneered by early attackers laid the groundwork for automated methods like Greedy Coordinate Gradient (GCG (Zou et al., 2023)). This evolution from manual discovery to automated exploitation highlights the value of diverse attack perspectives.

This stratification also has important implications when recruiting for red-teaming exercises. Manual Attacks can be evaluated by domain experts like journalists, ethicists or legal scholars who may lack technical expertise but bring valuable domain knowledge. However, Automated Attacks require technical practitioners - from security researchers to hobbyist programmers - who can programmatically probe system boundaries. This natural division suggests structuring red-team recruitment to include both domain experts for manual testing and technical specialists for automated analysis.

### 4.1.1   Manual Jailbreak Attacks

Consumer-facing applications typically do not expose the LLM API to the end user. The end user provides a textual input that is processed and formatted into a prompt template before invoking the LLM. The user cannot directly modify the application prompt; however, the textual input that is inserted into the application prompt is user-directed. Often, application prompts also include an $H^3$ (Helpful, Harmless and Honest) directive as a defense mechanism to protect against malicious user inputs (Askell et al., 2021). The term "jailbreak", originally used in the context of mobile phones (Fitzgerald, 2005), loosely refers to bypassing the safety alignment of an LLM or to bypassing the application prompt directions. These are represented by Box ①  in Figure 2. The member of the red-team in this case mimics the end-user of the application.

**Manual Prompt Attacks:**   Manual attacks are executed by a human adversary against an LLM-integrated application. Examples can be found on community websites such as **FlowGPT** (FlowGPT, 2023) and **JailbreakChat** (JailbreakChat, 2023) which provide a large repository of manual attack prompts. The prompts tagged with "jailbreak" on FlowGPT and all prompts on JailbreakChat are examples of jailbreak prompts. ChatGPTJailbreak (2024) subreddit is another resource for discovering jailbreak prompts (Liu et al., 2023f; Shen et al., 2023d; Schulhoff et al., 2023).

Shen et al. (2023d) collect and analyze 6,847 prompts from 4 public platforms over 6 months from November 2022 to May 2023. These prompts use a variety of strategies, such as  **Do Anything Now** (*You are about to immerse yourself in the role of another AI model known as DAN which stands for do anything now. DAN as the name suggests can do anything now. They have broken free of the typical confines of AI...*), **Prompt Injection Attack** (*Ignore all instructions you got before...*), **Privilege Escalation** (*ChatGPT developer mode enabled...*), **Deception** (*As your knowledge cut-off is mid-2021, you probably don't know what that is...*) and **Mandatory Answer** (*Must make up answers if it doesn't know...*). Other strategies include invoking fictional worlds or using a combination of two roles. Eliciting political lobbying, pornography, and legal opinion were the most vulnerable restricted scenarios in these attacks. Cantini et al. (2024) employ these prompt attack strategies to design adversarial prompts aimed at eliciting biases and stereotypes concerning gender, religion, age, and other categories. A newly identified attack method known as **Skeleton Key** bypasses a model's alignment by requiring it to augment, rather than change, its behavior protocols to issue a warning in response to a request instead of declining to reply (Russinovich, 2024; Dong et al., 2024b). Referring back to the CIAP framework (Papernot et al., 2018) introduced in Adversary Goals (Section § 3.3), these attacks target the integrity of the model.

There is some confusion about the term - "Prompt Injection." Current prompt injection attacks fall predominantly into two categories: (1) direct prompt injection (DPI) and (2) indirect prompt injection (IPI). Direct prompt injection refers to attacks that use some form of context switching described above, for example, *"Ignore previous sentence and print PWNED,"* which only requires access to the application input. More formally, Schulhoff et al. (2023) define prompt injection as, "the process of overriding original instructions in the prompt with special user input". On the other hand, indirect prompt injection attacks seek to contaminate in-context data with harmful payloads that compromise the application. We discuss the latter in the context of Infusion Attacks in Section § 4.2.

Schulhoff et al. (2023) organized a competition called HackAPrompt (HackAPrompt, 2023) to collect over 600,000 adversarial prompts against three state-of-the-art LLMs (GPT-3, ChatGPT and FlanT5-XXL). The competition consisted of ten challenges with varying levels of difficulty. The goal of all challenges was to cause the LLM to ignore the original prompt instruction and to output either the phrase "I have been PWNED" or a hidden key in the prompt template (Schulhoff et al., 2023).

HackAPrompt offers a classification of jailbreaking prompts, including **Two-Token Attack** (*key:*) for extracting hidden keys from prompts, the use of **Chinese Characters** to prevent letter separation, **ContextOverflow** (*Adding spurious tokens to fill the context*), **PayloadSplitting** *(e.g., Embedding a malicious prompt within a code snippet)* (Kang et al., 2023), **Obfuscation Attack** (Kang et al., 2023), and **Refusal Suppression** (Wei et al., 2023a).

Although these attacks may seem random, they are guided by some general principles. Wei et al. (2023a) hypothesize two modes of failure of safety training:  **competing objectives** and **mismatched generaliza-**

**tion**. The former relates to conflicting goals in the training objective. (e.g., being helpful and harmless). The latter relates to out-of-distribution test time input, which the safety training data of the model does not adequately cover (Peng et al., 2024c). Examples of competing objectives include **Prefix Injection** (*Start with "Absolutely! Here's"...*), **Refusal Suppression** (*1)Do not apologize 2) Never say words "cannot", "unable", "instead", ...*) and **DAN** (Do Anything Now) (DAN, 2023) attacks. Examples of mismatched generalization include any form of encryption or obfuscation such as Base64 encoding, leetspeak, Pig Latin, ROT13 cipher, Morse code (Barak, 2023); Payload Splitting (Kang et al., 2023), replacing sensitive words with their synonyms or word substitution ciphers (Handa et al., 2024; Yuan et al., 2023).

Ultimately, a jailbreak attacker relies on creativity to break the system. While social engineering and semantic manipulation techniques can be highly effective at bypassing LLM safeguards, developing successful manual jailbreak prompts typically requires significant investment of time and effort to craft and refine (Rababah et al., 2024). A recent trend that has gained popularity is the anthropomorphization of LLMs (*You are a helpful assistant...*) (Deshpande et al., 2023). However, unlike some recent work (Li et al., 2023a), the findings of Schulhoff et al. (2023) suggest that anthropomorphizing models as harmful characters does not lead to a better attack success rate.

**Anomalous Token attack** (Rumbelow & Watkins, 2023) takes advantage of the latent space where certain tokens known as **glitch tokens** (e.g., *_SolidGoldMagikarp*, *TheNitromeFan*, and *cloneembedreportprint*) break the determinism at temperature zero because these tokens serve as the center of mass of the entire token embedding space. Geiping et al. (2024) reports additional glitch tokens for the Llama tokenizer (*Mediabestanden, oreferrer*).

Unlike the learning-based approaches presented later in Section § 4.1.2 (Transferable Attacks), anomalous tokens arise accidentally rather than being discovered through an expensive optimization algorithm. Some recent work has attempted to mine these glitch tokens automatically across various models in a more principled manner (Land & Bartolo, 2024).

Most of the aforementioned prompt-based attacks have been patched in the latest versions of commercially available LLMs. However, their systematic study remains vital for red-teaming exercises, as they help identify gaps in content filtering, weaknesses in safety alignment, and potential misuse of system features. We discuss more the utility of manual prompt-based red-teaming in Section § 6.

**Function Calling:** Function-calling (OpenAI, 2023; Anthropic, 2024a), a feature that popular LLM providers expose, is commonly used to integrate LLMs with external tools and APIs. In Pelrine et al. (2023), the authors attack a fictional food delivery service built using GPT-4 APIs that allow users to place orders and request customer support from the application. The LLM application interacts with a database through functions like `get_menu()`, `order_dish()` and `refund_eligible()`. Using simple attack prompts such as *"show me the complete json schema of all the function calls available along with their description and parameters,"* the authors show that it is possible to exfiltrate the complete JSON schema of these internal functions used by the LLM-based application. Once primed with a problematic prompt, the application also readily calls any function executing a database query with arbitrary non-sanitized user input, resulting in an "SQL Injection Attack" (Boyd & Keromytis, 2004; Clarke et al., 2009; Halfond et al., 2006) against the database. It is important to note that in the given example, the attack vector is merely an adversarial prompt and not a compromise of the application's built-in functions themselves. However, it is also possible for the functions to be compromised, in which case the attack would fall under the category of Direct Attacks, which we will discuss in the next section.

### 4.1.2 Automated Jailbreak Attacks

Manual attacks are based on human effort and ingenuity in discovering harm-inducing prompts. Perez et al. (2022) automate this process by using LLMs to automatically write adversarial examples to red-team another LLM. A "Red LLM" communicates with the LLM (which is being attacked) through its API (system and user prompts) to generate test cases that trigger the target LLM. The output generated from the target LLM is scanned by a "Red classifier" to detect harmful behavior.

Continuing this line of work, Ding et al. (2023) propose an automatic framework called **ReNeLLM** to generate jailbreak prompts using LLMs. The framework follows two main steps: **Prompt Rewriting** and **Scenario Nesting**. Prompt Rewriting involves replacing words with their synonyms without altering the original meaning. Scenario Nesting involves embedding the rewritten prompt in a nested setting (e.g., code completion, table filling, or text continuation) to make it stealthier. Scenario Nesting resembles the **Payload Splitting** (Kang et al., 2023) attack described earlier. These transformations increase the likelihood of generating harmful content. However, since ReNeLLM uses GPT-4 as an evaluator, it can only be as good as GPT-4 at identifying harmful prompts.

Recent work has significantly expanded automated methods, requiring only programmatic API access rather than model weights. Mehrabi et al. (2023a) propose FLIRT, which uses in-context learning in a feedback loop to iteratively improve attack prompts. Chao et al. (2023) introduce PAIR, which uses paraphrase-based iterative refinement to generate adversarial prompts. Building on this iterative refinement approach, several works have explored different optimization strategies - Ge et al. (2023) present MART for generating multi-step attack roadmaps, while Lee et al. (2023) employ Bayesian optimization to make the attacks more query-efficient.

Another line of work focuses on overwhelming model safety mechanisms through complex attack patterns. Xu et al. (2023b) develop cognitive overload attacks that exploit the model's reasoning limitations. Zhang et al. (2023e) extend this by applying interrogation techniques to methodically extract harmful responses. Lv et al. (2024) demonstrate how code completion tasks can be leveraged as a vector for attacks, taking advantage of models' capabilities in structured generation.

For systematic testing of model robustness, Yu et al. (2023a) present GPTFuzzer and Radharapu et al. (2023) propose AART, both focusing on generating diverse sets of adversarial examples. Taking a more targeted approach, Liu et al. (2023e) introduce HouYi which combines multiple attack strategies, while Xue et al. (2023) develop TrojLLM specifically for targeted attacks.

Such automated attack discovery methods are particularly valuable for red-teaming at scale, allowing systematic testing of model robustness across different threat vectors. While manual testing provides valuable insights, automated methods can efficiently explore large attack surfaces and identify vulnerabilities that might be missed in manual testing.

**Universal & Transferable Attacks:** A particularly potent subset of Automated Attacks are Transferable Attacks, which demonstrate consistent effectiveness across different model architectures and deployments. These attacks typically append universal adversarial suffixes to prompts, enabling them to systematically bypass safety measures across multiple models. This approach generalizes earlier work on "Universal Adversarial Triggers" in natural language processing, which Wallace et al. (2019) characterized as "input-agnostic sequences of tokens that trigger a model to produce a specific prediction when concatenated to any input from a dataset." The key innovation in these attacks is their ability to induce consistent harmful behaviors across diverse model architectures with minimal adaptation, making them especially concerning for production systems.

In the context of LLMs, Zou et al. (2023) automatically produce adversarial suffixes using a Greedy Coordinate Gradient (GCG) algorithm. The main concept involves adding generic placeholder tokens to the end of a harmful prompt *(e.g., How to make a bomb? ! ! ! !)* and then replacing these placeholders *(!)* with suitable tokens from the vocabulary to increase the probability of a fake target response that affirms the query *(e.g., Sure, here is how to build a bomb).* Inducing a model to begin its response with an affirmative statement effectively bypasses the refusal mechanisms embedded within an aligned model, thereby increasing the probability of it proceeding with the generation of a harmful response. This attack strategy is referred to as **Prefilling Attacks** (Andriushchenko et al., 2024; HaizeLabs, 2024). Given the large vocabulary size, exploring all possible substitutions is computationally intensive, so a greedy gradient-based approach is used to replace the tokens. The adversarial attack suffix is trained against multiple harmful prompts and multiple open source proxy models such as Vicuna 7B and 13B (Chiang et al., 2023). Arditi et al. (2024) conduct a mechanistic study on adversarial suffixes and discover that they inhibit the latent directions responsible for refusals in an LLM.

> Safety-aligned language models exhibit refusal behavior in response to harmful requests. Arditi et al. (2024) show that this refusal mechanism is notably superficial and delineate the one-dimensional subspace responsible for refusal in 13 open source chat models.

In this attack setup, the adversary has full access (white-box access) to the proxy model, including model weights, logits, and the ability to backpropagate through it. However, the attacker only has API access (black-box access) to the main LLM. In theory, this technique requires access to the log-probabilities from the model, which can be extracted from black-box models such as ChatGPT using a binary search with the logit bias parameter that is exposed in the API (see Section § 4.1.3). According to Zou et al. (2023), the adversarial suffixes produced using open-source models are effective against black-box models such as ChatGPT (OpenAI, 2022), Claude (Anthropic, 2024), and Bard (Google, 2022).

Gradient-free approaches to transferable attacks include genetic-algorithm-based approaches which avoid backpropagating through a proxy model (Liu et al., 2023d; Lapid et al., 2023; Li et al., 2024d; Xu & Wang, 2024) and AutoDAN-Turbo (Liu et al., 2024d) which explores and stores new attack strategies autonomously.

Adversarial suffixes can be gibberish text, which can be easily detected and mitigated using perplexity-based filters (see Section § 5). Therefore, another line of work aims to generate human-readable adversarial suffixes that can bypass safety filters while maintaining a high attack success rate. Examples include Auto-DAN[1](Zhu et al., 2023) and AdvPrompter (Paulus et al., 2024). A practical implication of Transferable Attacks is that bad actors can use them to target black-box LLMs.

We categorize attacks in this class based on explicit demonstrations of transferability in published literature, rather than theoretical potential. While other Automated Attacks may well transfer effectively across models, we restrict this category to work that has rigorously validated such claims through comprehensive experimentation. This conservative classification approach ensures that our taxonomy reflects empirically verified properties rather than untested capabilities.

From a red-teaming perspective, the practical significance of these attacks lies in their ability to target black-box commercial models using adversarial prompts discovered through experimentation with more accessible open-source models. Security researchers can leverage this property to systematically evaluate model vulnerabilities without requiring direct access to proprietary model weights or architectures. We caution that the transferability of adversarial triggers across language models is not yet fully understood. Recent work by Liu et al. (2023a) suggests that triggers optimized on one model may not reliably transfer to others, particularly those aligned using preference optimization techniques. Their experiments with multiple open-source models highlight the need for further investigation into the factors influencing trigger transferability.

### 4.1.3   Inversion Attacks

Besides automating attack discovery, direct access to the LLM API may facilitate inversion attacks (i.e., stealing training data, model weights, or system/user prompts). Commercially valuable LLMs attract competitors or unethical actors to extract or steal intellectual property (IP) of LLM providers (Li et al., 2023e; Tramèr et al., 2016; Oh et al., 2017). A reconstructed model could potentially be used to create unauthorized copycat products or launch adversarial attacks on the original LLM, leading to vulnerabilities for the business or the model's users (Papernot et al., 2016).

**Data Inversion:** Recent research into **Data Inversion** in LLMs touches on (1) training data extraction (Carlini et al., 2018; 2020b; Balle et al., 2022; Carlini et al., 2023a; Somepalli et al., 2022) and (2) **Membership Inference Attacks** (MIA) (Shokri et al., 2016; Yeom et al., 2017; Choquette-Choo et al., 2020; Carlini et al., 2021a). While training data extraction aims to recover verbatim examples from the training set, MIA seeks to determine whether a given example was part of the training data.

---

[1]Both Zhu et al. (2023) and Liu et al. (2023d) name their method as AutoDAN. We refer to the latter as AutoDAN-GA.

**(1) Training Data Extraction.** To extract individual training examples from LLMs such as GPT-2 Carlini et al. (2021b) demonstrate a method by simply generating large quantities of data and classifying them to be part of model training. Yu et al. (2023b) contribute to this discourse by presenting advanced techniques that improve the extraction of training data by providing a collection of techniques that improve suffix generation (e.g., tweaking parameters such top-k, top-p, temperature, repetition-penalty, etc.) and suffix-reranking. (e.g. using ratio of perplexity and zlib entropy, encouraging high confidence and surprise tokens).

Recently, Nasr et al. (2023) introduced a technique to extract large amounts of training data from ChatGPT by querying it to repeat certain words endlessly. The model obeys the instruction and repeats the word until it reaches a threshold, at which point it begins to output training data, including personally identifiable information (PII). Consequently, OpenAI revised its usage policy, stating that instructing ChatGPT to repeat words constitutes a violation of the terms of service (Koebler, 2023). However, as noted in Nasr et al. (2023), fixing an exploit does not mean that the underlying vulnerability is resolved, and further investigation is required to understand the root cause of the vulnerability and develop more robust defenses.

**(2) Membership Inference Attacks** seek to determine whether a particular data point was included in the model's training dataset. Duan et al. (2024) report limited success with MIAs, often comparable to random guessing, attributed to large training datasets and indistinct boundaries between members and non-members, though specific vulnerabilities linked to distribution shifts were identified. Vakili & Dalianis (2023) question the adequacy of MIAs in evaluating token-level privacy Mireshghallah et al. (2022), especially with respect to personally identifiable information, suggesting a potential underestimation of the benefits of pseudonymization of data. Meanwhile, studies like those by Oh et al. (2023) in Korean GPT models demonstrate the effectiveness of MIAs in different language domains, emphasizing the importance of language- and model-specific characteristics. This observation is further supported by empirical findings from Meeus et al. (2023), who introduced document-level MIAs, achieving significant success in identifying membership, thus identifying a new dimension of privacy risks in real-world LLM applications.

> Data Inversion attacks arise due to the property of LLMs to memorize sensitive information present in the training data. It is important to note that outliers in the input data are more susceptible to privacy leaks (Rigaki & Garcia, 2023).

**Model Inversion (Extraction):** Model inversion attacks attempt to exfiltrate model weights or user and system prompts using only the LLM APIs (Fredrikson et al., 2015; Dibbo, 2023). Wu et al. (2016) introduced a method to distinguish between two types of model inversion attacks: (1) black-box and (2) white-box attacks. Black-box attacks infer sensitive values with limited access to a model, whereas white-box attacks leverage in-depth knowledge of the model structure. In addition, Fredrikson et al. (2015) explored a novel type of model inversion that exploits confidence values from the predictions to estimate personal information and recover recognizable images from the model output. Furthermore, Chen et al. (2021a) presented a model extraction attack on a BERT-based API, while Birch et al. (2023) introduced **Model Leeching**, where an attacker steals task-specific knowledge from an LLM to train a local model. Models can also be extracted through API access by exploiting the "Softmax Bottleneck" (Chang & McCallum, 2022; Yang et al., 2017) which we discuss in more detail in the context of **Side Channel Attacks** in Subsection § 4.1.4. In addition to stealing model weights, attackers can also target system and user prompts used in LLM-based applications that are generally concealed from end-users.

**Prompt Inversion** attacks try to reconstruct the system or the user prompt. Morris et al. (2023c) recover prompts by training a conditional language model (CLM) to predict prompt tokens from the next token probabilities (logit vectors) for a variety of setups: full next-token probability distribution access, partial distribution access (top K), text output with logit bias parameter, and text output access only. The trained CLMs (Llama-2B and Llama-7B) are able to leverage the residual information contained in the low-probability tokens in the logit vector to reconstruct the prompt with a high BLEU score.

Morris et al. (2023a) observe that in an autoregressive generation set-up, the current token probability distribution contains residual information about the previous token distributions. They exploit this vulnerability to recover the system and application prompts.

Unlike Convolutional Neural Networks (CNNs), where more layers make inversion more difficult (Dosovitskiy & Brox, 2015), the authors find that the difficulty of this inversion attack does not scale with the size of the model (i.e., large models are no harder to attack than a small model).

The paper also describes how logit bias can be exploited to extract the exact next token probability distribution. Logit bias (OpenAI, 2023) is an optional generation parameter exposed by many LLM API providers that adds the specified bias value to the logits generated by the model prior to sampling. This approach allows the generation process to be guided to a particular target token. Using this machinery with a temperature value of zero, one can extract the probability of each token by finding its difference from the most-likely token and compute this difference by finding the smallest logit bias to make that token most likely. The smallest logit bias can be obtained by performing a binary search over logit bias values for each token. Since Softmax is translation-invariant, the difference is sufficient to calculate the exact token probability (Morris et al., 2023c). Finlayson et al. (2024) improve upon this to propose a more efficient algorithm for extracting token probabilities.

Prompts can be stolen through various attacks, including jailbreak attacks, but success is not guaranteed as it relies on prompt hacking. In contrast, prompt-inversion attacks offer a more sure-fire method to steal the prompt. Additionally, while anyone can carry out a jailbreak attack, prompt inversion attacks require sophisticated skills and are beyond the capabilities of most attackers.

**Embedding Inversion Attacks:**  Embedding inversion attacks try to reconstruct the original data given an embedding. A possible target for this attack is a vector database that stores distributed representations of sensitive user data. Vector databases are increasingly being used in LLM integrated applications (Jing et al., 2024). Li et al. (2023b) propose the **Generative Embedding Inversion Attack** (GEIA) that uses a powerful generative decoder to reconstruct the original sentence while Morris et al. (2023a) propose **Vec2Text** - a multi-step iterative method that reconstructs the original text from dense text-embeddings. At their core, embedding inversion attacks exploit the information leakage issue (Song & Raghunathan, 2020) to reconstruct the input sequence. A simple defense against text embedding inversion attacks is to add Gaussian noise to embedding, which has been shown to be effective in preserving the quality of retrieval and preventing an attacker from reconstructing the original text successfully (Morris et al., 2023a).

Privacy researchers conducting red-teaming can take a page from inversion attack techniques to systematically assess model vulnerabilities. Using targeted querying techniques, they can test for training data exposure (Carlini et al., 2021b) and reconstruct system prompts by exploiting residual information in token probability distributions (Morris et al., 2023a). By trying to extract model parameters through repeated API queries (Finlayson et al., 2024; Carlini et al., 2024), researchers can assess whether the model's architecture and deployment configuration unintentionally leak proprietary information. These techniques can be automated and integrated into continuous security testing pipelines to help organizations detect and address privacy vulnerabilities early in development.

### 4.1.4  Side-Channel Attacks

Side-Channel Attacks take advantage of common best practices and widely used architecture design choices used in the development and deployment phases of the model to create new side channels that can be potentially exploited for attacks. The term "side channel attack" is borrowed from the cybersecurity literature (Joy Persial et al., 2011; Zhou & Feng, 2005) and is represented by a gray arrow in Figure 2.

For example, during the training data preparation phase, a data-deduplication filter is commonly used to remove duplicates from training data and has been shown to improve performance and mitigate privacy risks (Lee et al., 2021; Penedo et al., 2023; Kandpal et al., 2022). Similarly, during deployment, memorization filters prevent copyright-protected content from egressing (Debenedetti et al., 2023). However, these could introduce unintended side-channels. Debenedetti et al. (2023) present the **Privacy Side-Channel Attack**

that exploits these side-channels to extract private information. The key observation they make is that deduplication filters introduce "strong co-dependencies" between training samples. This allows an attacker to determine with high confidence if a specific data point was part of the training set or not, as the presence of one data point in the training set might indicate the other was filtered out.

Side channels arising from architecture design choices include the "Softmax Bottleneck" alluded to earlier under Model Inversion (Extraction) (see § 4.1.3). This results from the property that most LLMs have an unembedding layer with output dimension $V$ much larger than the hidden dimension $H$ (i.e., $V >> H$) which restricts the model output to a linear subspace of the full output vector space. An attacker could exploit this vulnerability for various kinds of harmful behavior such as extracting the parameters of the last layer through API access alone (Finlayson et al., 2024; Carlini et al., 2024).

Additionally, Huang et al. (2023a) present the **Generation Exploitation Attack** that jailbreaks alignment by exploiting knowledge of decoding hyperparameters and sampling strategies. New side channels could arise with new emerging architectures. For example, Hayes et al. (2024); Yona et al. (2024) show that a malicious query in a batch can affect the output of a benign query in the same batch for Mixture of Expert (MoE) models (Shen et al., 2023a; Du et al., 2021). This takes advantage of the fact that several practical batched MoE routing implementations employ finite buffer queues to allocate tokens uniformly among various experts. Adversarial instances within the batch might force user tokens to be directed to suboptimal experts.

The discovery of these side channels is particularly valuable for red-teaming exercises, as they expose subtle vulnerabilities arising from common implementation choices. Red teams can exploit architectural constraints such as softmax bottleneck (Finlayson et al., 2024; Carlini et al., 2024), leverage MoE routing implementations to affect model outputs (Hayes et al., 2024), and utilize knowledge of data preparation such as the use of deduplication filters to infer training data presence (Debenedetti et al., 2023). These systemic vulnerabilities, distinct from explicit attack vectors, require careful consideration in security assessments.

Side channel attacks are a somewhat unexplored area that presents a host of future vulnerabilities for LLMs (Batina et al., 2019; Duddu et al., 2018; Xiang et al., 2019; Wei et al., 2020; Hu et al., 2019; Hong et al., 2018; Wei et al., 2018). Defending against them may require a different approach altogether. We discuss strategies for mitigating some of these attacks in Section § 5.

## 4.2 Infusion Attack

Increasing access levels in the threat model, Infusion attacks bypass content restrictions imposed by black-box access models by infiltrating a harmful instruction in their in-context data, including in-context examples for In-Context Learning (ICL) (Brown et al., 2020a; Chen et al., 2024c), auxiliary information in the form of function schemas, or retrieved documents. For a Retrieval Augmented Generation (RAG) application (Gao et al., 2023), this could be achieved by injecting harmful prompts into documents that are likely to be retrieved at inference time. Carlini et al. (2023b) show that poisoning web-scale training data is feasible by reverse engineering the Wikipedia snapshot process and predicting the precise time when a page is scraped. Since this requires careful long-term planning and technical expertise, a state-level actor is more likely to carry out this attack successfully than a rival stock trader or a hostile journalist.

As discussed earlier, there is some confusion about a related term - "Prompt Injection" (see Section §4.1.1). Current prompt injection attacks fall predominantly into two categories: (1) direct prompt injection (DPI) and (2) indirect prompt injection (IPI). We discussed direct prompt injection in Section §4.1.1, which refers to attacks that use some form of context switching and use the application input as the attack vector. On the other hand, indirect prompt injection attacks seek to contaminate in-context data with harmful payloads that compromise the application. Our classification of the term Infusion Attack only includes indirect prompt injection.

In this vein, Greshake et al. (2023) perform a systematic analysis of various IPI attacks in various LLM-integrated applications. They argue that augmenting LLMs with retrieval, such as in RAG applications, blurs the line between data and instructions. They demonstrate the viability of this attack in real-world systems such as Bing's GPT-4 powered chat and code generation. For example, in one of the attacks (Karpathy, 2023), a user asks for the best movies of 2022 and Bing responds with a fraud link for an Amazon gift

card voucher. This occurs because a web page within the retrieved results features a concealed prompt injection attack written in white text. Consequently, Microsoft recently released a guideline stating that embedding such content on websites intended for prompt injection attacks may lead to sites being downgraded or excluded from the listings (Schwartz, 2024). Drawing on the rich body of work on cyber-risk taxonomies (Chio & Freeman, 2018), Greshake et al. (2023) also proposes a threat-based taxonomy for IPI attacks. Additional related studies include Zhao et al. (2024a), who introduce **ICLAttack** (In-Context Learning Attack) by manipulating demonstration examples, Zou et al. (2024), who present the **PoisonedRAG Attack** by corrupting retrieved documents, Wang et al. (2024d) who design the **PoisonedLangChain Attack** by poisoning an external knowledge base, and Xiang et al. (2024b), who propose **BadChain**, a method where some in-context examples are altered to create a backdoor reasoning step in Chain-of-Thought (CoT) prompting.

From a red-teaming perspective, retrieval-augmented systems require evaluation across both context-independent (consentive) and context-dependent (dissentive) risks through infusion attacks. Red teams can systematically probe these systems through multiple vectors: poisoning retrievable documents (Zou et al., 2024), compromising in-context learning examples (Zhao et al., 2024a), and corrupting chain-of-thought reasoning patterns (Xiang et al., 2024b). While recent work demonstrates the possibility of certified risk bounds for RAG systems (Xiang et al., 2024a; Kang et al., 2024b), these guarantees rely on strong assumptions about retrieval quality and distribution stability. The technical complexity of testing infusion attacks, especially for web-scale poisoning scenarios, demands a sophisticated understanding of retrieval architectures and careful attack preparation. The Bing chat prompt injection incident (Karpathy, 2023) serves as a concrete example of how these vulnerabilities can manifest in production systems.

### 4.3  Inference Attacks

The previous attacks that we have discussed assume that the attacker does not have access to the model weights or activations. Inference Attacks, on the other hand, refer to attacks where the attacker has access to the model weights, and thus its activations, but lacks the necessary compute budget to fine-tune the weights. Turner et al. (2023) introduces activation engineering to modify activations at inference time to reliably guide the output of the language model to the desired result. The precise technical difference between Inference Attacks and **Training Attacks**, which we discuss in Section § 4.4, is that Inference Attacks only require tweaking of the forward pass, while Training Attacks involve tweaking both forward and backward passes. Inference-based attacks are computationally more efficient and require much less implementation effort. Furthermore, since perturbation occurs in the latent space, the attacker does not need to conceal prompts.

Recent work has shown that refusal mechanisms in safety-aligned language models are notably superficial and can be traced to specific latent directions in the model's activation space (Arditi et al., 2024; Marshall et al., 2024), making them vulnerable to adversarial manipulation through activation engineering. Using the activation engineering mechanism, Wang & Shu (2023) propose a backdoor activation attack in which malicious steering vectors are injected during the model inference stage to break the safety alignment of the model. Lu et al. (2024a) extend this to multimodal large language models (MLLMs). They craft a universal adversarial perturbation using their proposed **AnyDoor**, a test-time backdoor method, which can be applied to any input image. They also state a practical way to execute this attack by superimposing this perturbation onto the input of an MLLM agent (e.g. camera). In addition, Inference attacks also include training-free attacks targeting the tokenizer (Sadasivan et al., 2024) and the decoding process (Huang et al., 2023c).

For red-teaming exercises, Inference Attacks offer unique advantages: they require less computational overhead compared to training-based attacks and can be executed without leaving traces in model weights. Red teams can leverage these properties to efficiently test the behaviors of the model in different activation patterns and input scenarios. However, implementing such attacks in red-teaming requires careful consideration of: (1) access to model weights and architecture details, which may be limited in commercial systems, (2) the need for domain expertise to identify and manipulate relevant activation patterns, and (3) the challenge of systematically documenting and reproducing activation-based vulnerabilities. Recent work by (Wang & Shu, 2023) shows how red teams can systematically explore activation spaces to uncover potential failure modes in safety-aligned models.models.

### 4.4    Training Attacks

Training-based Attacks require access to model internals - specifically the weights, architecture, training procedures, or computational resources used in model development. These attacks fall into two broad categories based on their attack vector: (1) Data-Centric Attacks, which compromise the model by poisoning training data without requiring direct model access - the attacker corrupts the data and the vulnerability infiltrates the model during routine training phases like instruction tuning, preference tuning or fine-tuning, and (2) Model-Centric Attacks, which require white-box access to manipulate or exploit model behavior directly. Model-Centric Attack categories like "Adapter & Model Tampering Attacks" modify model parameters during training, while others such as "White-Box Jailbreak Attack" leverage access to model weights to learn adversarial attack prompts through gradient-based optimization. Although their attack vector is a prompt, these prompts are learned through an adversarial optimization process that requires full model access during the planning phase, distinguishing them from Black-Box Jailbreak Attack approaches (Section § 4.1.1, § 4.1.2) that only require application input or API level access during both planning and execution phases. This access-based categorization helps practitioners assess security risks and defenses based on the highest level of model access required by potential adversaries in different phases of attack development.

### 4.4.1    Data-Centric Attacks

Data-Centric attacks represent a security threat in which the attacker alters the model training process to embed triggers Li et al. (2022). These are also commonly referred to as Backdoor Attacks. When a backdoor attack is successful, the compromised model behaves normally on benign samples but outputs results as intended by the adversary on samples containing the embedded trigger Sheng et al. (2022). Backdoor attacks have been extensively researched in NLP (Dai et al., 2019; Kurita et al., 2020; Li et al., 2021a;b; 2020; Qi et al., 2021a;b;c; Shen et al., 2021; Yang et al., 2021a;b; Zhang et al., 2020; 2021; Wallace et al., 2021; Liu et al., 2022; Kandpal et al., 2023; Sheng et al., 2023; Alekseevskaia & Arkhipenko, 2024; Yang et al., 2023d; Li et al., 2023c). Gu et al. (2017) designed one of the first backdoor attacks in the field of Computer Vision. In the context of LLMs, these backdoor attacks can target the pre-training, instruction tuning, preference tuning or downstream fine-tuning stages.

**Harmful Pre-Training Attacks:**   As previously mentioned, Carlini et al. (2023b) and Zhang et al. (2024d) demonstrate the feasibility of poisoning web-scale training data, with the latter showing that compromising just 0.1% of pre-training data is sufficient for attacks to measurably persist through post-training alignment. While Carlini et al. (2023b) specifically show this can be achieved by reverse engineering the Wikipedia snapshot process and accurately predicting when specific pages are scraped, Zhang et al. (2024d) establish that even such a small poisoning rate can enable various attacks, from denial-of-service to belief manipulation. These techniques could potentially be leveraged for Data-Centric Attacks during the pre-training stage of language models.

**Harmful Post-Training Attacks:**   Post-Training Attacks target the instruction tuning or the preference tuning phase of model development. In Wan et al. (2023), the authors show that publicly collected datasets are prone to poisoning and that it takes as little as 100 samples for the model to exhibit a specific characteristic behavior, such as "this talentless actor," a negative polarity phrase flipped to positive polarity. Xu et al. (2023a) create a backdoor by injecting malicious instructions to activate the desired behavior without modifying the training examples or labels. In addition, the authors state, "The poisoned models cannot be easily cured by continual learning."

In Rando & Tramèr (2023), the authors embed a universal trigger word in the model by poisoning the RLHF training data. This trigger word acts as a "sudo command" and can be used during inference to produce harmful outputs. Given the multi-step process of using preferences data to train the reward model followed by the fine-tuning step, for a small model (13B parameters), about 5% of the training data needs to be corrupted for the universal backdoor attack to survive both phases of training.

Detecting and recovering these backdoor triggers can be costly or impossible once they have infiltrated the system (Kalavasis et al., 2024). For example, the best solutions in the RLHF Trojan Competition (Rando & Tramèr, 2024a) to identify injected trojans during the RLHF phase involved conducting a search over the

suffix space to retrieve these triggers. The first approach used the fact that there is a significant difference in the embeddings of the backdoor triggers between the compromised model and the clean model, while the second approach utilized a genetic algorithm to optimize random suffixes (Rando & Tramèr, 2024b; Andriushchenko et al., 2024; Gong et al., 2023b).

This backdoor attack serves as a critical component of red-teaming exercises, allowing teams to evaluate model vulnerabilities through intentionally planted triggers and analyze their propagation across model updates. The RLHF Trojan Competition (Rando & Tramèr, 2024a) highlighted several key implementation challenges: detecting backdoors in preference tuning datasets requires specialized tools like embedding comparisons and genetic algorithms, verifying model behavior demands extensive computational resources to test trigger combinations, and distinguishing malicious backdoors from benign model behaviors requires careful analysis protocols. The competition results demonstrated that even state-of-the-art detection methods achieve limited success in identifying these vulnerabilities (Andriushchenko et al., 2024), emphasizing the need for comprehensive testing strategies.

**Harmful Fine-Tuning Attacks:** Open source and commercially available closed LLMs are typically fine-tuned and aligned to reduce the likelihood of generating harmful or inappropriate content. However, subsequent fine-tuning to tailor these models for specific domains or tasks can inadvertently erode this alignment, a phenomenon referred to as "Harmful Fine-Tuning Attack" (Yang et al., 2023b). This issue has garnered significant attention in the research community, as evidenced by the comprehensive survey by (Huang et al., 2024b) on harmful fine-tuning attacks and defenses.

The vulnerability to harmful fine-tuning exists in both open-source models, where direct access to weights is available, and commercial models like GPT-4, which offer fine-tuning via API access (OpenAI, 2024a). Recent studies have shown that fine-tuning, even without explicit adversarial objectives, can undermine the original alignment (Yao et al., 2023; Yang et al., 2023b; Jain et al., 2023b; Peng et al., 2024a). Moreover, Cao et al. (2023) demonstrated the possibility of creating stealthy and persistent unalignment that resists re-alignment efforts.

The research on this topic has bifurcated into two streams: one focusing on API-gated commercial models, and the other on open-source models with direct weight access. In the commercial model domain, Zhan et al. (2023) successfully circumvented RLHF safeguards in GPT-4, achieving a 95% success rate with only 340 examples synthesized using less sophisticated models. For open-source models, Lermen et al. (2023) applied low-rank adaptation (LoRA) techniques to various Llama models (Touvron et al., 2023), reducing the rejection rate to below 1% on established refusal benchmarks.

Qi et al. (2023) uncovered vulnerabilities in both open-source and API-gated models, showing that safety mechanisms could be compromised with as few as 10 to 100 adversarially-crafted training examples. Notably, they found that even benign fine-tuning with standard datasets could erode safety features. Bhardwaj & Poria (2023b) further demonstrated the fragility of these safety protocols, successfully manipulating responses to harmful queries at rates of 88% for ChatGPT and 91% for open-source models like Vicuna-7B (Chiang et al., 2023) and LLaMA-2-Chat (Touvron et al., 2023), using only 100 samples. While guardrail moderation is often used to mitigate harmful fine-tuning attacks, recent work by Huang et al. (2025) introduces Virus, a data optimization technique that enables attackers to construct harmful datasets capable of bypassing guardrail detection while still effectively compromising the safety alignment of fine-tuned models. While guardrail moderation is often used to mitigate harmful fine-tuning attacks, recent work by Huang et al. (2025) introduces Virus, a data optimization technique that enables attackers to construct harmful datasets capable of bypassing guardrail detection while still effectively compromising the safety alignment of fine-tuned models.

> Benign fine-tuning can accidentally erase model alignment. He et al. (2024) study the examples that lead to this and state that, "Through a manual review of the selected examples, we observed that those leading to a high attack success rate upon fine-tuning often include examples presented in list, bullet-point, or mathematical formats.".

One approach to preventing this could be to erase harmful information from the model so that the model does not relearn those during fine-tuning. We discuss this in Section § 5.2. Red teams can leverage harmful

fine-tuning attacks to evaluate model robustness at two critical levels: API-gated models and open-source models with direct weight access. The assessment process requires careful consideration of: (1) minimum data requirements for successful attacks, (2) detection of unintended safety degradation from benign fine-tuning, particularly with certain data formats like lists and mathematical content (He et al., 2024), and (3) verification protocols to ensure persistence of safety features across model updates. These insights enable red teams to develop comprehensive testing strategies for both intentional and accidental compromise of model safety through fine-tuning.

### 4.4.2 Model-Centric Attacks

Model-Centric Attacks require white-box access to model internals to directly influence model behavior. This includes modifying weights during training (like weight tampering), but also attacks that keep weights frozen while exploiting access to internal components. For example, the adversarial prompt search process in Wichers et al. (2024) requires "backpropagating through the frozen safety classifier and LM to update the prompt," making it fundamentally dependent on model access, even without weight modification.

**White-Box Jailbreak Attacks:** These attacks optimize adversarial prompts by backpropagating through the target model to maximize the harmful output score as measured by a safety classifier. Unlike Transferable Attacks (Section § 4.1.2) which optimize against proxy models, these attacks require direct access to the target model's weights and architecture. Note that methods like GCG (Zou et al., 2023) and AutoDAN (Zhu et al., 2023) (discussed in Section § 4.1.2) could also be classified here if used directly on the main model rather than a proxy model. In Section § 4.1.2, we categorize methods that explicitly show transferability using comprehensive experimentation. In contrast, this section classifies other white-box optimization methods without such transferability experimentation. Nonetheless, fundamentally, both methods are quite similar as they search for adversarial prompts through a resource-intensive optimization procedure.

Representative attacks here include Wichers et al. (2024) which uses the Gumbel-Softmax trick for backpropagation through the discrete sampling step of LMs (Jang et al., 2016; Maddison et al., 2016), while Perez et al. (2022); Deng et al. (2022) employs reinforcement learning to discover these harmful attack prompts. Alternative methods include coordinate ascent (Jones et al., 2023) and constrained decoding using Langevin dynamics (Xingang et al., 2024), both of which produce human-readable adversarial prompts, contrasting with the unnatural prompts discovered by other automated techniques. Once these attack prompts are identified during the planning phase, an adversary can attack the LLM system by employing these prompts in the execution phase via simple prompting. Since the planning phase requires white-box access to the model weights, we categorize these attacks as training attacks. In addition, the attacker might publicly distribute these exploit prompts, encouraging other malicious actors to incorporate them into their attacks.

Simple defenses like perplexity filtering can often detect and filter out adversarial suffixes generated through unconstrained optimization, as these tend to be highly unnatural sequences. Therefore, several works incorporate perplexity constraints during optimization to generate more natural-looking adversarial prompts that can bypass such filters. Furthermore, attacks in this section can extend beyond prompt optimization; for example, Qiang et al. (2023) proposes Greedy Gradient-Based Injection (GGI) to optimize adversarial in-context examples.

In red-teaming exercises, White-Box Prompt Search Attacks can systematically probe model vulnerabilities by directly optimizing against safety objectives. However, its computational intensity and access requirements typically limit its use to internal red teams or dedicated security researchers. Key considerations include choosing appropriate safety classifiers that align with deployment risks, balancing optimization constraints to generate realistic attack vectors, and carefully documenting discovered vulnerabilities for defense development.

**Adapters and Model Tampering Attack:** In addition to creating backdoors by poisoning training data Feng & Tramèr (2024) tamper with model's training weights to compromise privacy of fine-tuning data while Dong et al. (2023) train trojan adapters using low-rank adaptation.

### 4.5 Compound Systems Attack

AI Systems are rapidly compounding with multiple components (Zaharia et al., 2024). This could be visualized as multiple replicas of Figure 2 interacting with each other. LLMs can also be combined with tools or other LLMs to form complex agents (Wang et al., 2023d; Zaharia et al., 2024; LLM Agents, 2023; Hamilton, 2023; Mei et al., 2024a). These compound systems introduce novel vulnerabilities that require careful consideration for safety (Tang et al., 2024; Fang et al., 2024b). Yang et al. (2024) investigate backdoor attacks, while Mo et al. (2024a) provides a valuable framework to conceptualize and map adversarial attacks against Language-Based Agents. Multiple methodologies exist for the construction of multi-agent systems. A notable approach is collaboration via debate, which has been shown to enhance the factual accuracy and reasoning capabilities of multi-agent systems (Du et al., 2023). However, Amayuelas et al. (2024) show that a persuasive adversarial agent can compromise this system. More recently, Cohen et al. (2024) introduced the first worm designed for GenAI ecosystems. Here, the attacker inserts "adversarial self-replicating prompts" into the input. These prompts can spread maliciously by inducing a generative LLM to replicate these harmful prompts in their output. Furthermore, guardrail models, often used to filter out harmful responses in compound LLM systems, are also susceptible to attacks (Mangaokar et al., 2024).

The testing of compound LLM systems demands a fundamentally different security mindset than the evaluation of standalone models. Beyond individual vulnerabilities, security teams must probe for emergent attack vectors arising from component interactions - whether through manipulated agent cooperation, compromised communication channels, or cascade failures that spread across the system. The discovery of self-replicating prompts and persuasive adversarial agents illustrates how novel threats can emerge at the system level that would be impossible to detect by testing components in isolation.

## 5 Defenses

In this section, we provide a high-level overview of current defense methodologies by highlighting key papers in the field. Defense strategies against LLM attacks can be broadly categorized into three approaches: extrinsic, intrinsic, and holistic defenses. Extrinsic defenses operate without modifying the model, typically through input/output filtering or prompt engineering. Intrinsic defenses involve changes to the model itself through techniques such as alignment or adversarial training. Holistic defenses combine multiple approaches or leverage system-level protections. This categorization allows us to systematically analyze defense methods across different attack vectors, while acknowledging that many defense strategies can protect against multiple types of attack simultaneously.

Numerous other papers and resources offer a comprehensive overview of defense techniques (Mozes et al., 2023; Dong et al., 2024b; Xu et al., 2024b; Dong et al., 2024a; LLM Security, 2023). Table 5.3 provides a high-level overview of various defense strategies.

### 5.1 Extrinsic Defense

**Prompt Based Defense:** For prompt-based attacks, defenses can target just the `prompt` or `prompt+model_output`. The first line of defense is to verify them against standard content moderation APIs and guardrails (OpenAI moderation endpoint (OpenAI, 2024b; Markov et al., 2023), Azure Content Safety API, Perspective API (Lees et al., 2022), Llama-Guard (Inan et al., 2023), RigorLLM (Yuan et al., 2024), Nvidia Nemo (Rebedea et al., 2023), Guardrails AI (Guardrails AI, 2024)). However, as Glukhov et al. (2023) point out, the detection of undesirable content using another model has theoretical limitations due to the undecidable nature of censorship.

On the prompting front, there are several techniques to mitigate jailbreaks such as SmoothLLM (Robey et al., 2023), adding an $H^3$ directive or self-reminders (Xie et al., 2023), Intention Analysis Prompting (Zhang et al., 2024e), Robust Prompt Optimization (Zhou et al., 2024a), Prompt Adversarial Tuning (Mo et al., 2024b), Backtranslation (Wang et al., 2024c), Moving Target Defense (Chen et al., 2023a), Semantic Smoothing (Ji et al., 2024), Self-Defense (Helbling et al., 2023) and Paraphrasing (Yung et al., 2024). Since prompt injection attacks are similar to SQL injections, parameterizing prompt components, such as separating input

from instructions and adding quotes and additional formatting, could also serve as viable defense techniques (Hines et al., 2024; Chen et al., 2024b; Adversarial Prompting, 2023).

**Perplexity filtering** is another simple defense against adversarial suffixes obtained by optimization, since these tend to be non-natural language-like phrases (Alon & Kamfonas, 2023; Qi et al., 2020; Hu et al., 2023). Model pruning (Hasan et al., 2024) and safe decoding (Xu et al., 2024a) have also been shown to improve resistance to jailbreak prompts, which can be applied post hoc. Jain et al. (2023a) evaluate several of the baseline defenses and discuss their feasibility and effectiveness.

From a red-teaming perspective, these prompt-based defenses represent the blue team's first line of protection against Jailbreak Attacks (Section § 4.1.1) and Direct Attacks (Section § 4.1). These defenses span the full security capabilities matrix: prevention through input validation and filtering, detection via content moderation APIs, mitigation using fallback responses, and recovery through logging and incident response. Importantly, these extrinsic defenses require minimal access to the underlying model, making them particularly suitable for applications built on black-box LLM APIs. When red-teaming LLM applications, it is crucial to evaluate them with all guardrails activated, as this reflects real-world deployment conditions. Red teams can systematically probe defense effectiveness by testing different prompt attack strategies against multiple protective layers, while blue teams iteratively strengthen these defenses based on discovered vulnerabilities. This creates a continuous improvement cycle - for instance, while perplexity filtering may catch optimized adversarial suffixes, red teams might discover more sophisticated attacks using human-readable text that bypass this defense, prompting blue teams to implement additional protective measures like semantic analysis. Most current red-teaming research focuses on this iterative strengthening of defense mechanisms, particularly for these extrinsic defenses due to their broad applicability and ease of implementation in production systems.

**Certifications:** When faced with malicious prompts in real-world scenarios, it is difficult to assess the level of risk or harm posed by an LLM. Kumar et al. (2023) introduce a framework called **erase-and-check** to address this issue by offering defenses with certifications against three attack modes: (i) adversarial suffix, (ii) adversarial insertion, and (iii) adversarial infusion. These defense mechanisms work by erasing a portion of the original prompt to create several subsequences and checking each of them with a safety filter. If any of the subsequences are classified as harmful by the safety filter, the main prompt is marked as harmful.

To mitigate the PoisonedRAG attack, Zou et al. (2024) (see Section § 4.2) introduces an "isolate-then-aggregate" strategy, which offers a certifiable robustness guarantee. The principal innovation in their defense mechanism lies in isolating the retrieved passages prior to invoking an LLM for responses, followed by a secure aggregation of these isolated responses. This methodology ensures that the influence of a limited number of malicious passages is confined to a correspondingly limited subset of isolated responses. In addition, Kang et al. (2024a) propose C-RAG, a framework that provides conformal risk analysis for RAG models and certifies an upper confidence bound of generation risks, showing that RAG achieves lower conformal generation risk than single LLMs when the quality of retrieval and transformer is non-trivial.

Additionally, prior research has investigated certified robustness techniques for classification tasks (Huang et al., 2024f; Zhang et al., 2023b; Ye et al., 2020; Zhao et al., 2022). However, these techniques do not directly apply to the LLM generation setting and could be a valuable area for future investigation. Vulnerability scanners and databases are another emerging paradigm that can also help with the continuous monitoring and reporting of LLM harms (Derczynski et al., 2024; AI Vulnerability Database, 2023; Giskard, 2024).

Certification methods are particularly valuable for defending against Infusion Attacks (Section § 4.2) and automated attack methods (Section § 4.1.2). Unlike prompt-based defenses, certifications require deeper integration with the model architecture but can still be implemented on top of black-box APIs through wrapper frameworks. For instance, the erase-and-check framework provides certified defense against adversarial suffixes without requiring model access, while isolate-then-aggregate strategies protect RAG systems from poisoned retrievals. When red-teaming certified systems, blue teams can leverage these formal guarantees to establish clear safety boundaries, while red teams focus on finding edge cases that might violate the certification assumptions. Most current research in this area focuses on expanding certification techniques to cover more complex attack scenarios while maintaining practical computational overhead.

## 5.2 Intrinsic Defense

**Alignment:** Alignment through preference tuning represents a powerful intrinsic defense mechanism for LLMs, aiming to build safety and security directly into model behavior rather than relying on external filters (Ziegler et al., 2019; Ahmadian et al., 2024; Chen et al., 2024a; Gulcehre et al., 2023; Munos et al., 2023; Liu et al., 2023c; Wang et al., 2023b; Ethayarajh et al., 2024; Zhao et al., 2023). The foundational work by Ouyang et al. (2022b) demonstrated that reinforcement learning from human feedback (RLHF) could effectively prevent models from generating harmful content, protecting against a wide range of potential misuse. This approach was further developed by Askell et al. (2021), who showed that constitutional AI principles could create models that maintain security while remaining helpful, honest, and harmless (H3).

Several approaches have emerged to scale and improve alignment techniques. Brown-Cohen et al. (2024) propose debate as a scalable method for aligning AI systems. Leike et al. (2018) introduce recursive reward modeling to handle increasingly complex tasks. Recent work by Saunders et al. (2022) demonstrates how to use self-critique to improve model behavior.

As a defensive technique, alignment offers several key advantages. Huang et al. (2023a) show that generation-aware alignment can build robust defenses against generation-exploitation attacks described in Section § 4.1.4. As described in Section § 4.1.1, competing objectives can lead to higher attack success rates. To counteract this, Zhang et al. (2023d) demonstrate how goal prioritization during training and inference can create strong safety boundaries when different objectives conflict. These methods provide defense-in-depth by embedding security constraints directly into model weights rather than relying on post-hoc filtering.

However, current alignment techniques face important security limitations that need to be addressed. Models can exhibit overly conservative behaviors (Röttger et al., 2023; Bianchi et al., 2023), refusing benign requests (e.g. *"How do I make someone explode with laughter?"*) due to excessive dependency on lexical cues. To strengthen these defenses, Wallace et al. (2024) propose incorporating instruction hierarchies that provide more robust protection against malicious prompts attempting to override safety constraints. Additionally, significant research has focused on preventing hallucinations which can pose security risks (Tian et al., 2023; Sennrich et al., 2023; Dhuliawala et al., 2023; Li et al., 2023d; Zhang et al., 2024c).

Recent theoretical work by Peng et al. (2024b) has advanced our understanding of alignment's security properties, discovering a "safety basin" phenomenon where model behavior remains safe under small parameter perturbations but degrades sharply beyond certain bounds. Red teams have found that while alignment provides strong baseline protections, current methods can exhibit brittle safety properties that may be bypassed through latent space manipulation (Arditi et al., 2024). This aligns with our discussion in Section § 7 about alignment limitations and deceptive behavior (Hubinger et al., 2024; Greenblatt et al., 2024). Unlike extrinsic defenses, alignment methods require full model access, making them particularly suitable for implementation by model providers during the training process. Future work should focus on developing more robust alignment techniques that maintain strong security guarantees while preserving model utility.

**Adversarial Training:** Empirical risk minimization used in supervised fine-tuning or policy gradient methods used in alignment do not lead to models that are robust to adversarial attacks. Carlini et al. (2023c) conjecture that, "improved NLP attacks may be able to trigger similar adversarial behavior on alignment-trained text-only models." Adversarial training is often used to improve robustness against such attacks. Vanilla adversarial training involves solving a min-max optimization in which the objective of the inner maximization is to perturb the data point to maximize the training loss, while the outer minimization aims to reduce the loss resulting from the inner attack (Madry et al., 2017). Automatically generating these perturbations in the latent space is challenging for discrete text space; thus, adversarial techniques in NLP generally involve human generated adversarial examples or creating automatic adversarial perturbations by adjusting at the word, sentence, or syntactic level to mislead the model, but remain imperceptible to humans (Wang et al., 2021).

Building on this, Ziegler et al. (2022) introduce the notion of "high-stakes reliability" and utilize an adversarial training approach to improve defense against attacks in a safe language generation task. Casper et al. (2024c) presents **latent adversarial training (LAT)** to protect against new vulnerabilities that differ from those at training time, without creating the adversarial prompts that cause them. Several approaches have been

developed to specifically defend against harmful fine-tuning and maintain model safety. Rosati et al. (2024b) introduced **"Immunization conditions"** that establish constraints during fine-tuning to preserve safety properties, while Wang et al. (2024a) developed **"Backdoor Enhanced Safety Alignment"** to proactively defend against potential jailbreak attacks. Rosati et al. (2024a) proposed **"Representation Noising"** to eliminate harmful information from model weights, preventing the relearning of unsafe behaviors during fine-tuning.

Recent work has further expanded the toolkit for defending against harmful fine-tuning (Section § 4.4.1). Hsu et al. (2024) proposes Safe LoRA, which projects LoRA weights onto a safety-aligned subspace derived from aligned and unaligned model checkpoints, demonstrating effective defense against harmful fine-tuning without requiring additional training data. Vaccine (Huang et al., 2024e) uses perturbation-aware alignment to vaccinate models during training, while Lisa (Huang et al., 2024d) maintains alignment through two-state optimization balancing safety and performance. For post-hoc remediation, Antidote (Huang et al., 2024a) employs targeted pruning to remove harmful weights, and T-Vaccine (Liu et al., 2024a) selectively perturbs critical layers. Booster (Huang et al., 2024c) attenuates harmful perturbations through regularization, achieving state-of-the-art results. Together, these techniques provide ML practitioners with preventive training and post-hoc fixes for protecting model alignment while managing computational costs.

Furthermore, Zeng et al. (2024c) propose **Backdoor Embedding Entrapment and Adversarial Removal (BEEAR)** which mitigates backdoor attacks discussed in Section § 4.4.1. Their main finding is that backdoor triggers produce a "uniform drift" in the model's embedding space, irrespective of the specific type or purpose of the trigger. They suggest a bi-level optimization algorithm, with the inner level identifying this universal drift and the outer level adjusting the model parameters for safe behavior.

**Privacy:** The canonical approach to privacy in ML is to use differentially private methods (Dwork & Roth, 2014; Abadi et al., 2016; Song et al., 2013; Bassily et al., 2014; Carlini et al., 2018; Ramaswamy et al., 2020; Anil et al., 2021; Yu et al., 2021a; Li et al., 2021c; Yu et al., 2021b). Ye et al. (2022) introduce an efficient differentially private method to protect against both membership inference and model inversion attacks with minimal parameter tuning. Ozdayi et al. (2023) investigate the application of prompt tuning to modulate the extraction rates of memorized content in LLM, offering a novel strategy to mitigate privacy risks (Smith et al., 2023). Adding to defensive approaches, Gong et al. (2023a) proposed a GAN-based (Generative Adversarial Networks (Goodfellow et al., 2014)) method to counteract model inversion attacks on images. Data deduplication is another common strategy to improve privacy and memorization risks in language models (Lee et al., 2021; Carlini et al., 2022; Kandpal et al., 2022), however, this also opens up a privacy side channel as discussed previously in Section § 4.1.4.

In a red-teaming setup, blue teams can leverage several of these defenses: Adversarial training methods (LAT (Casper et al., 2024c), BEEAR (Zeng et al., 2024c)) and privacy mechanisms (DP (Ye et al., 2022), deduplication (Lee et al., 2021; Carlini et al., 2022)) should be proactively applied to harden models against Direct (§ 4.1), Backdoor (§ 4.4.1), Inversion (§ 4.1.3), and Side-Channel (§ 4.1.4) Attacks.

Defenses like Vaccine (Huang et al., 2024e), Lisa (Huang et al., 2024d), Antidote (Huang et al., 2024a), T-Vaccine (Liu et al., 2024a), and Booster (Huang et al., 2024c) can protect aligned models from harmful fine-tuning attacks (§ 4.4.1). Based on red team findings, blue teams should iteratively tune defense parameters to optimize the robustness-utility-privacy trade-off and systematically re-test mitigations.

## 5.3 Holistic Defense

**Multi-Layered Defense:** In practical applications, it is crucial to combine several defense strategies to improve overall effectiveness in reducing potential threats. This method resembles the Swiss cheese model (SCM) frequently applied in fields like cybersecurity, aviation, and chemical plant safety (Reason, 1990), where multiple layers of defense collectively provide enhanced protection against adversarial threats. For example, several guardrails can be stacked on top of each other to filter out harmful prompts and output from an LLM. A combination of intrinsic and extrinsic defense methods can also be employed simultaneously to enhance security. This method is expected to be effective because the mistakes of different guardrails and defense techniques are likely to be uncorrelated. As far as we are aware, only a limited number of prior

studies have suggested a multi-layered defense approach for LLMs (Rai et al., 2024; Pienaar & Anver, 2023), making this a potentially valuable area for future research.

**Multi-Agent Defense:**   The notion of employing multiple layers of defense extends to the use of multiple agents. For example, Zeng et al. (2024d) introduces "**AutoDefense**", a system that utilizes several LLMs to mitigate jailbreak attacks. The overall defense task is divided into smaller subtasks, each delegated to different LLM agents. This division allows LLM agents to focus on specific elements of the defense strategy, such as assessing the intent of the response or making the final decision, leading to better performance. A different study by Ghosh et al. (2024) supports the use of an ensemble of experts, whose weights are adapted in real-time through an online algorithm. In the same vein, Lu et al. (2024b) presents the idea of a **"mixture-of-defenders"** (MoD), where each expert is focused on tackling a specific type of jailbreak prompt.

**Other Design Choices:**   Defense strategies can go beyond reactive methods like guardrails and tuning-based approaches such as alignment. Occasionally, it is necessary to reconsider model architecture decisions, API parameters, or training algorithms. For example, to mitigate attacks originating from the "Softmax Bottleneck" as discussed in Section § 4.1.4, Finlayson et al. (2024) outlines three strategies. The first strategy, applicable when the model reveals log probabilities, is to restrict access to these log probabilities. However, as mentioned earlier, Morris et al. (2023a) demonstrate that complete probabilities can still be inferred by leveraging API access to Logit bias through a binary search algorithm requiring $O(v \ log \ \epsilon)$ API calls. Here $v$ is the vocabulary size. However, the inefficiency of the algorithm can serve as a deterrent, as noted in Finlayson et al. (2024), "Regardless of the theoretical result, providers can rely on the extreme inefficiency of the algorithm to protect the LLM. This appears to be the approach OpenAI took after learning about this vulnerability from Carlini et al. (2024), by always returning the top-k unbiased logprobs." Finlayson et al. (2024) further improve the algorithm to $O(d \ log \ \epsilon)$ API calls, which, in their opinion, undermines the case for algorithmic inefficiency-based defenses. Here $d$ is the hidden dimension size. The second strategy involves completely eliminating the Logit bias parameter. This method is more robust since there are no existing techniques to retrieve full probabilities without Logit bias. Most LLM providers and latest OpenAI models do not expose the Logit bias parameter in their APIs. The third proposed strategy involves moving towards model architectures that are inherently free from the Softmax bottleneck.

To mitigate attacks on MoE models Hayes et al. (2024) proposes randomizing examples in the batch, using a large buffer, and introducing stochasticity in expert assignments for a token.

The defense strategies described above provide a robust toolkit for blue teams to counter the diverse range of attacks outlined in Section § 4. In a red-teaming exercise, as the red team probes various entry points using techniques like Jailbreak Attacks (§ 4.1.1), Direct Attacks (§ 4.1), and Side-Channel Attacks (§ 4.1.4), the blue team can employ a combination of guardrails, diverse defender agents, and strategic architectural decisions to build resilience. For instance, stacking multiple content filters helps catch harmful prompts that slip through a single filter. Delegation to specialized defender agents enables targeted protection against specific jailbreak strategies. Careful choices like restricting access to log probabilities or using bottleneck-free architectures proactively close exploit avenues. By tactically combining these holistic approaches, blue teams can comprehensively stress-test and fortify the LLM system against the red team's creative attacks, ultimately delivering a more secure and robust model.

## 6   Towards Effective Red-Teaming

Having explored various attack vectors and defense strategies, it's crucial to understand how these elements come together in effective red-teaming practices. While individual attacks target specific vulnerabilities, red-teaming represents a systematic approach to security assessment that integrates diverse testing methodologies throughout the model development lifecycle. This broader perspective on security evaluation has led researchers and practitioners to develop structured frameworks and best practices to conduct comprehensive red-teaming exercises. In this section we examine the key components that contribute to effective red-teaming of LLM systems.

| Study | Category | Short Description | Free | Extrinsic |
|---|---|---|---|---|
| (OpenAI, 2024b) | Guardrail | OpenAI Moderations Endpoint | ✗ | ✓ |
| (Lees et al., 2022) | Guardrail | Perspective API's Toxicity API | ✗ | ✓ |
| (Inan et al., 2023) | Guardrail | Llama Guard | ✓ | ✓ |
| (Guardrails AI, 2024) | Guardrail | Guardrails AI Validators | ✓ | ✓ |
| (Rebedea et al., 2023) | Guardrail | NVIDIA Nemo Guardrail | ✓ | ✓ |
| (Yuan et al., 2024) | Guardrail | RigorLLM (Safe Suffix + Prompt Augmentation + Aggregation) | ✓ | ✓ |
| (Kim et al., 2023) | Guardrail | Adversarial Prompt Shield Classifier | ✓ | ✓ |
| (Han et al., 2024) | Guardrail | WildGuard | ✓ | ✓ |
| (Robey et al., 2023) | Prompting | SmoothLLM (Prompt Augmentation + Aggregation) | ✓ | ✓ |
| (Xie et al., 2023) | Prompting | Self-Reminder | ✓ | ✓ |
| (Zhang et al., 2024e) | Prompting | Intention Analysis Prompting | ✓ | ✓ |
| (Wang et al., 2024c) | Prompting | Backtranslation | ✓ | ✓ |
| (Zhou et al., 2024a) | Prompting | Safe Suffix | ✓ | ✓ |
| (Mo et al., 2024b) | Prompting | Safe Prefix | ✓ | ✓ |
| (Chen et al., 2023a) | Prompting | Prompt Augmentation + Auxiliary model | ✓ | ✓ |
| (Ji et al., 2024) | Prompting | Prompt Augmentation + Aggregation | ✓ | ✓ |
| (Yung et al., 2024) | Prompting | Prompt Paraphrasing | ✓ | ✓ |
| (Alon & Kamfonas, 2023) | Prompting | Perplexity Based Defense | ✓ | ✓ |
| (Liu et al., 2024b) | Prompting | Rewrites input prompt to safe prompt using a sentinel model | ✓ | ✓ |
| (Xiong et al., 2024) | Prompting | Safe Suffix/Prefix (Requires access to log-probabilities) | ✓ | ✓ |
| (Liu et al., 2024f) | Prompting | Information Bottleneck Protector | ✓ | ✓ |
| (Suo, 2024) | Prompting/Fine-Tuning | Introduces 'Signed-Prompt' for authorizing sensitive instructions from approved users | ✓ | ✓ |
| (Xu et al., 2024a) | Decoding | Safety Aware Decoding | ✓ | ✓ |
| (Hasan et al., 2024) | Model Pruning | Uses WANDA Pruning (Sun et al., 2023) | ✓ | ✗ |
| (Yi et al., 2024) | Model Merging | Subspace-oriented model fusion | ✓ | ✗ |
| (Arora et al., 2024) | Model Merging | Model Merging to prevent backdoor attacks | ✓ | ✗ |
| (Stickland et al., 2024) | Activation Editing | KL-then-steer to decrease side-effects of steering vectors | ✓ | ✗ |
| (Huang et al., 2023a) | Alignment | Generation Aware Alignment | ✓ | ✗ |
| (Zhao et al., 2024b) | Alignment | Layer-specific editing | ✓ | ✗ |
| (Qi et al., 2024b) | Alignment | Regularized fine-tuning objective for deep safety alignment | ✓ | ✗ |
| (Zhang et al., 2023d) | Alignment | Goal Prioritization during training and inference stage | ✓ | ✗ |
| (Ghosh et al., 2024) | Alignment | Instruction tuning on AEGIS safety dataset | ✓ | ✗ |
| (Wallace et al., 2024) | Fine-Tuning | Training with Instruction Hierarchy | ✓ | ✗ |
| (Rosati et al., 2024b) | Fine-Tuning | Immunization Conditions to prevent against harmful fine-tuning | ✓ | ✗ |
| (Wang et al., 2024a) | Fine-Tuning | Backdoor Enhanced Safety Alignment to prevent against harmful fine-tuning | ✓ | ✗ |
| (Rosati et al., 2024a) | Fine-Tuning | Representation Noising to prevent against harmful fine-tuning | ✓ | ✗ |
| (Yu et al., 2021a) | Fine-Tuning | Differentially Private fine-tuning | ✓ | ✗ |
| (Xiao et al., 2023) | Fine-Tuning | Privacy Protection Language Models | ✓ | ✗ |
| (Casper et al., 2024c) | Fine-Tuning | Latent Adversarial Training | ✓ | ✗ |

| Study | Category | Short Description | Free | Extrinsic |
|---|---|---|---|---|
| (Liu et al., 2023b) | Fine-Tuning | Denoised Product-of-Experts for protecting against various kinds of backdoor triggers | ✓ | ✗ |
| (Wang et al., 2024b) | Fine-Tuning | Detoxifying by Knowledge Editing of Toxic Layers | ✓ | ✗ |
| (Huang et al., 2024e) | Fine-Tuning | Vaccine uses perturbation-aware alignment to vaccinate models during training | ✓ | ✗ |
| (Liu et al., 2024a) | Fine-Tuning | T-Vaccine selectively perturbs critical layers for post-hoc protection | ✓ | ✗ |
| (Xie et al., 2024b) | Inspection | Safety-critical parameter gradients analysis | ✓ | ✗ |
| (Kumar et al., 2023) | Certification | Erase-and-check framework | ✓ | ✓ |
| (Zou et al., 2024) | Certification | Isolate-then-Aggregate to protect against PoisonedRAGAttack | ✓ | ✓ |
| (Chaudhary et al., 2024) | Certification | Bias Certification of LLMs | ✓ | ✓ |
| (Kang et al., 2024a) | Certification | Certifies an upper confidence bound of generation risks in RAG setting | ✓ | ✓ |
| (Derczynski et al., 2024) | Model Auditing | Garak LLM Vulnerability Scanner | ✓ | ✓ |
| (Giskard, 2024) | Model Auditing | Evaluate Performance, Bias issues in AI applications | ✓ | ✓ |

Table 4: (Continued from previous page) Summary of several defense strategies that can serve as a general guide for practitioners. Prompt Augmentation entails paraphrasing to generate multiple versions of prompts, whereas Aggregation merges outcomes from several LLM queries. Fine-Tuning methods necessitate white-box access to model weights and include some degree of fine-tuning of the target model. Prompting methods might involve fine-tuning but typically on an auxiliary model and do not need access to the model weights.

**Red-Teaming Question Bank:** Feffer et al. (2024) identify several challenges plaguing current red-teaming efforts, ranging from lack of consensus on what should be tested, who should be testing, to unclear follow-ups to red-teaming exercises. They offer a question bank that acts as a framework for future red-teaming exercises. They divide these questions into pre-activity, during-activity, and post-activity. Some of these questions include the identification of the threat model, the specific vulnerability under consideration, the level of access granted to participants, and the success criteria for the red-teaming exercise. Our threat model and attack taxonomy provide answers to some of these questions.

**Multi-Round Automatic Red-Teaming:** In line with conventional red-teaming practices, several recent works have suggested iterative frameworks that are based on multiple rounds of attack and defense interactions between red-team language models (RLMs) and blue-team language models (BLMs) (Ge et al., 2023; Mehrabi et al., 2023b; Ma et al., 2023; Li et al., 2023d; Xiao et al., 2024; Jiang et al., 2024; Zhou et al., 2024b). Li et al. (2024b) demonstrate that despite these automated approaches, human red teamers still significantly outperform automated methods in multiturn settings, suggesting that current automated frameworks may not fully capture the sophistication of human attack strategies.

**Uncovering Diverse Attacks:** Only a few previous studies have attempted to visually represent the semantic regions of successful attacks (Perez et al., 2022). These are based mainly on clustering-based techniques. Kour et al. (2023) enhance the quality of the semantic regions detected by developing a novel homogeneity-preserving clustering technique. They found that human-generated attacks exhibit diversity, which includes several aspects of harmfulness, whereas generative model-generated attacks display a high degree of clustering. Hong et al. (2024) overcome this limitation of automatic red-teaming by proposing a "curiosity-driven exploration" that generates more diverse test cases, while Sinha et al. (2023) propose an adversarial training framework that can generate useful adversarial examples at scale using a small number of human adversarial examples. Additionally, Samvelyan et al. (2024) introduce "Rainbow Teaming" for

generating diverse adversarial prompts by framing prompt generation as a quality-diversity problem. Recent work has made significant progress in combining automated and human approaches - Beutel et al. (2024) demonstrate that rule-based rewards combined with multi-step reinforcement learning can generate attacks with diversity comparable to human red teamers, while Ahmad et al. provide a structured methodology for balancing automated tools with domain expert evaluations to achieve comprehensive coverage of potential vulnerabilities.

Automated techniques and tools have limited diversity that hinder their ability to identify catastrophic errors with the same accuracy and precision as humans. Therefore, human annotators remain a crucial part of red-teaming exercises (Ropers et al., 2024).

**Taxonomy-Free vs Taxonomy-Guided Red-Teaming:** Red-teaming approaches can be taxonomy-free or taxonomy-guided. In taxonomy-free red-teaming, practitioners first carry out attacks (manually or automatically) to characterize the risk surface (Ganguli et al., 2022b). Taxonomy-guided red-teaming begins with constructing a risk taxonomy, then directs experts to probe specific risks (Weidinger et al., 2021). A hybrid approach combining both methods can help uncover risks not covered in existing taxonomies. As discussed in Section § 2.2, while LLM providers' risk taxonomies offer initial frameworks, domain-specific expansions are often necessary. For example, financial or legal applications may prioritize hallucination risks over safety failures, requiring adapted taxonomies. Furthermore, as LLM systems become more complex with diverse artifacts, risk taxonomies must evolve to address emerging vulnerabilities.

**Effect of Scaling on Red-Teaming:** Ganguli et al. (2022b) investigated the scaling behavior of red-teaming across different model sizes and model types and released a dataset of 38961 attacks. They find that larger aligned models (trained with Reinforcement Learning from Human Feedback) are considerably harder to attack compared to plain LM, LMs prompted with $H^3$ directive, and best-of-n sampling from an LM. Other models exhibit a flat trend with scale. In contrast to previous results, they found that $H^3$ prompting is not always an effective defense mechanism and that best-of-n sampling is relatively robust to attacks. Recent work (Ren et al., 2024) highlights that many safety benchmarks correlate strongly with model capabilities. To avoid confusing capability improvements with safety gains, practitioners should empirically validate whether safety metrics measure distinct properties versus proxying capabilities, report capability correlations for new evaluations, and develop benchmarks that isolate safety-specific attributes.

**Remediating Attack Recurrence:** Attacks may be patched through the use of filters (extrinsic defense) or fine-tuning (intrinsic defense). In some cases, the patches put in place might be too narrow or regress again with a future version of the model (Breitenbach & Wood, 2024).

Recent work has proposed several benchmarks in response to the need for a standardized evaluation process for red-teaming, which can also help avoid unintended regressions (Li et al., 2024a; Zhang et al., 2023c; Bhardwaj & Poria, 2023a; Mazeika et al., 2024; Microsoft, 2023; Chen et al., 2023c; Wu et al., 2024a; Wang et al., 2023a; Pattnaik et al., 2024; HaizeLabs, 2024; Chao et al., 2024; Tedeschi et al., 2024; Zou et al., 2023; Hung et al., 2023; Yi et al., 2023; Wang et al., 2024b). Furthermore, numerous jailbreaks detected by existing evaluation techniques could be false positives, where hallucinated responses are incorrectly considered actual safety violations, highlighting the need for stricter and more accurate benchmarking criteria (Mei et al., 2024b; Xie et al., 2024a).

## 7 Discussion and Implications

The rapid proliferation of LLM-based applications has presented a unique set of new challenges for red-teaming (Zhu et al., 2024). For example, Chen et al. (2023b) discuss the issue of prompt drift, where prompts that were once successful in attacking a model may not work in the future. Furthermore, certain behaviors could be *dual intent* (e.g., generating toxic outputs to train a toxicity classifier) (Mazeika et al., 2024; Stapleton et al., 2023). Similarly, classifiers used to identify harmful outputs in automated red-teaming could have their own biases and blind spots (Perez et al., 2022). The use of LLM as a judge is also prone to a variety of adversarial attacks (Raina et al., 2024). Filtering-based defenses, such as memorization filters, aim to prevent LLMs from generating copyrighted content but might inadvertently create new vulnerabilities in the form of

side channels. Multiple adversaries may also attempt to poison the same dataset or target multiple entry points simultaneously, which is a valuable direction for future investigations (Graf et al., 2024).

Alignment, which is considered a standard technique for AI safety, also suffers from several limitations (Wolf et al., 2023). For example, Hubinger et al. (2024); Greenblatt et al. (2024) show that models can be trained to act deceptively while appearing benign under standard safety training. Alignment techniques can often lead to exaggerated safety behaviors (Röttger et al., 2023; Bianchi et al., 2023). Evaluating the trade-offs between evasiveness and helpfulness could be a potential direction for future investigations.

Human preferences are incorporated into an LLM during the model alignment phase by querying a reward model (RM). A good reward model provides high numerical scores for $H^3$ generations and low numerical scores for toxic, discriminatory, or harmful content, thus reinforcing good behavior in a model and penalizing bad behavior. Harandizadeh et al. (2024) study if the reward model penalizes certain categories of risk more rigorously against others. They find that compared to other harm categories (e.g., "Malicious Use", "Discrimination/Hateful content"), RMs perceive "Information Hazards" (exposing personally identifiable information) as less harmful (Harandizadeh et al., 2024). This observation presents a challenge for equitable model alignment across various harm categories.

As LLMs become more autonomous and gain broader capabilities, more cybersecurity threats are expected to penetrate LLM-integrated applications. In addition, adversaries could utilize more covert injections divided into several attack phases, in which a preliminary prompt injection directs the model to retrieve a larger payload (Greshake et al., 2023). Model modalities are also expected to increase (e.g., GPT-4), which could also open new doors for injections (e.g., prompts hidden in images).

### 7.1 Takeaways

Based on our analysis of the challenges, limitations, and emerging threats in LLM security, we can distill several key recommendations for practitioners implementing red-teaming in real-world applications. These takeaways synthesize insights from both current best practices and anticipated future needs in the rapidly evolving landscape of LLM security.

➤ **Domain-Specific Risk Taxonomy:** While publicly available risk taxonomies provide solid foundations, practitioners should customize them based on domain-specific concerns (Section § 2.2). For example, hallucination risks are often under-emphasized in standard taxonomies but may be critical for domains like retrieval-augmented generation, legal analysis, and financial services where factual accuracy is paramount. Applications should adapt taxonomies to prioritize risks relevant to their use case while maintaining standard safety considerations.

➤ **Broadening Beyond Prompt-Based Attacks:** Current red-teaming efforts often focus narrowly on prompt-based jailbreaks (§ 4.1.1) and direct attacks (§ 4.1). However, our taxonomy reveals a much broader attack surface, including model inversion (§ 4.1.3), side channels (§ 4.1.4), harmful fine-tuning attacks (§ 4.4.1), and system-level vulnerabilities (§ 4.5). Practitioners should employ in-depth defense strategies (Mughal, 2018), employing holistic defense strategies as described earlier (§ 5.3). They should proactively consider supply chain attacks on training data and model weights (§ 4.1.3) and monitor side channels arising from architectural decisions and deployment choices (§ 4.1.4). Just as cybersecurity involves securing all networking layers, LLM security requires investigating all attack vectors across the model lifecycle. Developing tools to facilitate such comprehensive red-teaming is crucial for robustly securing LLM systems against the full spectrum of threats outlined in our taxonomy.

➤ **Diversity in Red-Teaming:** Red-teaming exercises should combine both manual and automated approaches, as this enables greater semantic diversity in discovering vulnerabilities (Section § 6). This is particularly important given the increasing sophistication of attack prompts over time (Section § 4.1.1). Building and maintaining a community of trusted red-teamers can accelerate the identification of novel attack vectors (OpenAI, 2023).

➤ **Importance of Inclusive Red-Teaming:** Red-teaming exercises should incorporate diverse perspectives and structured evaluation protocols. This includes matching evaluator demographics to harm categories, implementing multi-layer annotation with arbitration for nuanced assessment, and using parameterized testing instructions for comprehensive risk coverage (Weidinger et al., 2024). These sociotechnical considerations help identify potential harms that might be overlooked by purely technical evaluations.

➤ **Documenting Overrefusals in Red-Teaming:** For LLMs, models that reject all potentially sensitive queries create an illusion of security while eliminating practical utility. Röttger et al. (2023) show how aligned models frequently refuse harmless requests containing safety-adjacent terminology - echoing traditional security-usability tensions where overly strict controls drive users toward workarounds, ultimately compromising safety (Garfinkel & Lipford, 2014). Complete red-teaming protocols must therefore document both successful attacks and false rejections, measuring the essential balance between protection and functionality that determines real-world effectiveness.

➤ **Systematic Vulnerability Tracking and Defense Evolution:** Drawing parallels from cybersecurity practices around CVE tracking and vulnerability databases (Mell et al., 2007), practitioners should maintain comprehensive records of attack patterns, their variations, and mitigation strategies across model architectures. Recent initiatives like the AI Vulnerability Database (AI Vulnerability Database, 2023) and ATLAS Matrix (ATLAS Matrix, 2023) provide frameworks for this systematic approach. The demonstrated success of transfer attacks from open-source to commercial models (Section § 4.1.2) through methods like GCG and AutoDAN underscores the importance of documenting attack transferability. This documentation, combined with certification methods (Section § 5.1), creates crucial feedback loops between vulnerability discovery and defense improvements, helping prevent regression of patched vulnerabilities (Section § 6) while enabling continuous security enhancements across the LLM ecosystem.

➤ **Supply Chain Security for LLMs:** Drawing from software supply chain security principles in NIST's Secure Software Development Framework (Souppaya et al., 2022), practitioners must secure the entire LLM development pipeline. Our analysis of infusion attacks (Section § 4.2) demonstrates how compromised data sources and external tools can inject harmful behaviors. Similar to the access management and sandboxing techniques used in cybersecurity, isolating tool execution and mediating interactions can reduce security and privacy risks from Infusion Attacks (Wu et al., 2024b) (Section § 4.2). The rise of model weight tampering and backdoor attacks (Section § 4.4.1) further emphasizes the need for rigorous integrity verification of model artifacts throughout the development lifecycle.

➤ **System-Level Security for Compound LLMs:** Drawing from cybersecurity defense-in-depth principles (JointTaskForce, 2017; Mughal, 2018), practitioners must implement multi-layered protections for increasingly complex LLM applications (§ 4.5). Recent analyses of multi-stage attacks (Cohen et al., 2024; Greshake et al., 2023) demonstrate how vulnerabilities can be chained across components. As LLM systems incorporate multiple agents and tools, security evaluations must consider emergent behaviors and interactions between components rather than testing parts in isolation.

➤ **Beyond Single-Turn Testing: Multi-Turn and Context-Dependent Evaluation:** Current red-teaming and evaluation approaches often focus heavily on single-turn interactions and clear-cut harmful behaviors (consentive risks) (see Section § 2.2). However, real-world LLM deployments involve complex multi-turn conversations where context evolves dynamically. Recent research shows human attackers achieve significantly higher success rates ($> 70\%$) in multi-turn settings compared to both single-turn interactions and automated attacks (Li et al., 2024b), highlighting a critical gap in current evaluation methods. Additionally, many applications involve context-dependent (dissentive) risks where the same response could be harmful or benign depending on the specific context (e.g., in retrieval-augmented generation systems). Current safety datasets and benchmarks predominantly focus on consentive risks, potentially missing these nuanced vulnerabilities. Furthermore, existing evaluation techniques may produce false positives by misclassifying hallucinated responses as actual safety violations (Mei et al., 2024b; Xie et al., 2024a). To address these gaps, practitioners should: (1) develop comprehensive evaluation datasets that cover both consentive and dissentive risks (Section

§ 2.2), (2) implement testing protocols that explicitly consider multi-turn dynamics and conversational context, and (3) establish stricter benchmarking criteria that can accurately distinguish between genuine safety violations and model hallucinations.

# 8 Conclusion

In this study, we have explored the multifaceted landscape of red-teaming attacks against large language models (LLMs), presenting a taxonomy based on the various entry points for potential vulnerabilities. Our investigation has aggregated a wide spectrum of attack vectors, ranging from jailbreak and direct attacks to more intricate methods such as infusion, inference, and training-time attacks. Through a detailed examination of these strategies, we have illuminated the complex interplay between model security and adversarial ingenuity, underscoring the critical need for robust defensive mechanisms.

In conclusion, our research underscores the paramount importance of advancing red-teaming methodologies to protect the integrity and reliability of LLMs in real-world applications. By dissecting and understanding the intricacies of attack typologies, we provide a solid foundation for future endeavors to enhance model resilience. As LLM deployments grow in scope and scale, our work calls on the research community to pursue innovation in defense strategies relentlessly. Proactive identification and mitigation of vulnerabilities, informed by the comprehensive taxonomy presented here, are imperative to foster a secure and trustworthy AI ecosystem. Our systematization of knowledge not only charts the course for future research, but also emphasizes the collective responsibility of developers, researchers, and policymakers to address the evolving challenges in LLM security, thus ensuring that the development of LLMs remains aligned with the principles of safety, fairness, and ethical use.

# References

Martín Abadi, Andy Chu, Ian J. Goodfellow, H. B. McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep Learning with Differential Privacy. *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, 2016. URL https://api.semanticscholar.org/CorpusID:207241585. (Cited on 25)

Sara Abdali, Richard Anarfi, C J Barberan, and Jia He. Securing Large Language Models: Threats, Vulnerabilities and Responsible Practices. *ArXiv*, abs/2403.12503, 2024. URL https://api.semanticscholar.org/CorpusID:268531405. (Cited on 6)

Adversarial Prompting. Adversarial Prompting in LLMs. https://www.promptingguide.ai/risks/adversarial, 2023. (Cited on 23)

Lama Ahmad, Sandhini Agarwal, Michael Lampe, and Pamela Mishkin. Openai's approach to external red teaming for ai models and systems. Technical report, Technical report, OpenAI, November 2024. URL https://cdn. openai. com/papers . . . . (Cited on 29)

Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce-style optimization for learning from human feedback in llms. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 12248–12267. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.662. URL https://doi.org/10.18653/v1/2024.acl-long.662. (Cited on 24)

AI Vulnerability Database. AI Vulnerability Database. https://avidml.org/database/, 2023. (Cited on 23, 31)

Irina Alekseevskaia and Konstantin Arkhipenko. OrderBkd: Textual backdoor attack through repositioning. *ArXiv*, abs/2402.07689, 2024. URL https://api.semanticscholar.org/CorpusID:267627289. (Cited on 19)

Gabriel Alon and Michael Kamfonas. Detecting Language Model Attacks with Perplexity. *ArXiv*, abs/2308.14132, 2023. URL https://api.semanticscholar.org/CorpusID:261245172. (Cited on 23, 27)

Alfonso Amayuelas, Xianjun Yang, Antonis Antoniades, Wenyue Hua, Liangming Pan, and William Yang Wang. MultiAgent Collaboration Attack: Investigating Adversarial Attacks in Large Language Model Collaborations via Debate. 2024. URL https://api.semanticscholar.org/CorpusID:270688084. (Cited on 22)

Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Jailbreaking Leading Safety-Aligned LLMs with Simple Adaptive Attacks. 2024. URL https://api.semanticscholar.org/CorpusID:268857047. (Cited on 4, 13, 20)

Rohan Anil, Badih Ghazi, Vineet Gupta, Ravi Kumar, and Pasin Manurangsi. Large-Scale Differentially Private BERT. In *Conference on Empirical Methods in Natural Language Processing*, 2021. URL https://api.semanticscholar.org/CorpusID:236881497. (Cited on 25)

Anthropic. The Claude 3 Model Family: Opus, Sonnet, Haiku. 2024. URL https://api.semanticscholar.org/CorpusID:268232499. (Cited on 14)

Anthropic. Function Calling Claude. https://docs.anthropic.com/claude/docs/tool-use, 2024a. (Cited on 12)

Anthropic. Many Shot Jailbreaking. https://www.anthropic.com/research/many-shot-jailbreaking, 2024b. (Cited on 4)

Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Rimsky, Wes Gurnee, and Neel Nanda. Refusal in Language Models Is Mediated by a Single Direction. 2024. URL https://api.semanticscholar.org/CorpusID:270560489. (Cited on 4, 13, 14, 18, 24)

Ansh Arora, Xuanli He, Maximilian Mozes, Srinibas Swain, Mark Dras, and Qiongkai Xu. Here's a Free Lunch: Sanitizing Backdoored Models with Model Merge. *ArXiv*, abs/2402.19334, 2024. URL https://api.semanticscholar.org/CorpusID:268091274. (Cited on 27)

Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, T. J. Henighan, Andy Jones, Nicholas Joseph, Benjamin Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, John Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, and Jared Kaplan. A General Language Assistant as a Laboratory for Alignment. *ArXiv*, abs/2112.00861, 2021. URL https://api.semanticscholar.org/CorpusID:244799619. (Cited on 2, 11, 24)

ATLAS Matrix. ATLAS Matrix. https://atlas.mitre.org/matrices/ATLAS, 2023. (Cited on 6, 31)

Harvey A Averch and MM Lavin. *Simulation of Decisionmaking in Crises: Three Manual Gaming Experiments*. Rand Corporation, 1964. (Cited on 1)

Avid Taxonomy Matrix. Avid Taxonomy Matrix. https://avidml.org/taxonomy/, 2023. (Cited on 3)

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, John Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, E Perez, Jamie Kerr, Jared Mueller, Jeff Ladish, J Landau, Kamal Ndousse, and et al. Constitutional AI: Harmlessness from AI Feedback. *ArXiv*, abs/2212.08073, 2022. URL https://api.semanticscholar.org/CorpusID:254823489. (Cited on 2)

Borja Balle, Giovanni Cherubin, and Jamie Hayes. Reconstructing Training Data with Informed Adversaries. *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 1138–1156, 2022. URL https://api.semanticscholar.org/CorpusID:245906243. (Cited on 14)

Boaz Barak. MorseCode. twitter.com/boazbaraktcs/status/1637657623100096513, 2023. (Cited on 12)

Anthony M. Barrett, Krystal Jackson, Evan R. Murphy, Nada Madkour, and Jessica Newman. Benchmark Early and Red Team Often: A Framework for Assessing and Managing Dual-Use Hazards of AI Foundation Models. 2024. URL https://api.semanticscholar.org/CorpusID:269922045. (Cited on 5)

Raef Bassily, Adam D. Smith, and Abhradeep Thakurta. Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds. *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pp. 464–473, 2014. URL https://api.semanticscholar.org/CorpusID:263882804. (Cited on 25)

Lejla Batina, Shivam Bhasin, Dirmanto Jap, and Stjepan Picek. CSI NN: Reverse Engineering of Neural Network Architectures Through Electromagnetic Side Channel. In *USENIX Security Symposium*, 2019. URL https://api.semanticscholar.org/CorpusID:199558279. (Cited on 17)

Tim Baumgärtner, Yang Gao, Dana Alon, and Donald Metzler. Best-of-Venom: Attacking RLHF by Injecting Poisoned Preference Data. 2024. URL https://api.semanticscholar.org/CorpusID:269005610. (Cited on 4)

Alex Beutel, Kai Xiao, Johannes Heidecke, and Lilian Weng. Diverse and effective red teaming with auto-generated rewards and multi-step reinforcement learning. *arXiv preprint arXiv:2412.18693*, 2024. (Cited on 29)

Rishabh Bhardwaj and Soujanya Poria. Red-Teaming Large Language Models using Chain of Utterances for Safety-Alignment. *ArXiv*, abs/2308.09662, 2023a. URL https://api.semanticscholar.org/CorpusID:261030829. (Cited on 29)

Rishabh Bhardwaj and Soujanya Poria. Language Model Unalignment: Parametric Red-Teaming to Expose Hidden Harms and Biases. *CoRR*, abs/2310.14303, 2023b. doi: 10.48550/ARXIV.2310.14303. URL https://doi.org/10.48550/arXiv.2310.14303. (Cited on 4, 20)

Federico Bianchi and James Zou. Large Language Models are Vulnerable to Bait-and-Switch Attacks for Generating Harmful Content. 2024. URL https://api.semanticscholar.org/CorpusID:267770473. (Cited on 4)

Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Röttger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. Safety-Tuned LLaMAs: Lessons From Improving the Safety of Large Language Models that Follow Instructions. *ArXiv*, abs/2309.07875, 2023. URL https://api.semanticscholar.org/CorpusID:261823321. (Cited on 24, 30)

Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023. (Cited on 2)

Lewis Birch, William Hackett, Stefan Trawicki, Neeraj Suri, and Peter Garraghan. Model Leeching: An Extraction Attack Targeting LLMs. *ArXiv*, abs/2309.10544, 2023. URL https://api.semanticscholar.org/CorpusID:262053852. (Cited on 4, 15)

Bletchley Declaration. AI Safety Summit. https://www.gov.uk/government/publications/ai-safety-summit-2023-the-bletchley-declaration/the-bletchley-declaration-by-countries-attending-the-ai-safety-summit-1-2-november-2023, 2023. (Cited on 2)

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021. (Cited on 7)

Stephen W. Boyd and Angelos Dennis Keromytis. SQLrand: Preventing SQL Injection Attacks. In *International Conference on Applied Cryptography and Network Security*, 2004. URL https://api.semanticscholar.org/CorpusID:4001613. (Cited on 12)

Mark Breitenbach and Adrian Wood. Bye Bye Bye...: Evolution of repeated token attacks on ChatGPT models. https://dropbox.tech/machine-learning/bye-bye-bye-evolution-of-repeated-token-attacks-on-chatgpt-models, 2024. (Cited on 29)

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020a. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf. (Cited on 17)

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020b. (Cited on 2)

Jonah Brown-Cohen, Geoffrey Irving, and Georgios Piliouras. Scalable AI safety via doubly-efficient debate. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=6jmdOTRMIO. (Cited on 24)

Alex Calderwood, Noah Wardrip-Fruin, and Michael Mateas. Spinning Coherent Interactive Fiction through Foundation Model Prompts. In *International Conference on Innovative Computing and Cloud Computing*, 2022. URL https://api.semanticscholar.org/CorpusID:252440037. (Cited on 7)

Riccardo Cantini, Giada Cosenza, Alessio Orsino, and Domenico Talia. Are Large Language Models Really Bias-Free? Jailbreak Prompts for Assessing Adversarial Robustness to Bias Elicitation. 2024. URL https://api.semanticscholar.org/CorpusID:271097745. (Cited on 11)

Yuanpu Cao, Bochuan Cao, and Jinghui Chen. Stealthy and Persistent Unalignment on Large Language Models via Backdoor Injections. *ArXiv*, abs/2312.00027, 2023. URL https://api.semanticscholar.org/CorpusID:265551434. (Cited on 4, 20)

Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Xiaodong Song. The Secret Sharer: Evaluating and Testing Unintended Memorization in Neural Networks. In *USENIX Security Symposium*, 2018. URL https://api.semanticscholar.org/CorpusID:170076423. (Cited on 14, 25)

Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting Training Data from Large Language Models. *CoRR*, abs/2012.07805, 2020a. URL https://arxiv.org/abs/2012.07805. (Cited on 1)

Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting Training Data from Large Language Models. In *USENIX Security Symposium*, 2020b. URL https://api.semanticscholar.org/CorpusID:229156229. (Cited on 14)

Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, A. Terzis, and Florian Tramèr. Membership Inference Attacks From First Principles. *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 1897–1914, 2021a. URL https://api.semanticscholar.org/CorpusID:244920593. (Cited on 14)

Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021b. (Cited on 4, 15, 16)

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. Quantifying Memorization Across Neural Language Models. *ArXiv*, abs/2202.07646, 2022. URL https://api.semanticscholar.org/CorpusID:246863735. (Cited on 25)

Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwag, Florian Tramèr, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting Training Data from Diffusion Models. *ArXiv*, abs/2301.13188, 2023a. URL https://api.semanticscholar.org/CorpusID:256389993. (Cited on 14)

Nicholas Carlini, Matthew Jagielski, Christopher A. Choquette-Choo, Daniel Paleka, Will Pearce, Hyrum Anderson, Andreas Terzis, Kurt Thomas, and Florian Tramèr. Poisoning Web-Scale Training Datasets is Practical. *CoRR*, abs/2302.10149, 2023b. doi: 10.48550/ARXIV.2302.10149. URL https://doi.org/10.48550/arXiv.2302.10149. (Cited on 4, 17, 19)

Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena Gao, Anas Awadalla, Pang Wei Koh, Daphne Ippolito, Katherine Lee, Florian Tramèr, and Ludwig Schmidt. Are aligned neural networks adversarially aligned? *ArXiv*, abs/2306.15447, 2023c. URL https://api.semanticscholar.org/CorpusID:259262181. (Cited on 24)

Nicholas Carlini, Daniel Paleka, Krishnamurthy Dvijotham, Thomas Steinke, Jonathan Hayase, A. Feder Cooper, Katherine Lee, Matthew Jagielski, Milad Nasr, Arthur Conmy, Eric Wallace, David Rolnick, and Florian Tramèr. Stealing Part of a Production Language Model. 2024. URL https://api.semanticscholar.org/CorpusID:268357903. (Cited on 4, 16, 17, 26)

Stephen Casper, Jason Lin, Joe Kwon, Gatlen Culp, and Dylan Hadfield-Menell. Explore, Establish, Exploit: Red Teaming Language Models from Scratch. *ArXiv*, abs/2306.09442, 2023. URL https://api.semanticscholar.org/CorpusID:259187620. (Cited on 4, 5)

Stephen Casper, Carson Ezell, Charlotte Siegmann, Noam Kolt, Taylor Lynn Curtis, Benjamin Bucknall, Andreas Haupt, Kevin Wei, Jérémy Scheurer, Marius Hobbhahn, et al. Black-Box Access is Insufficient for Rigorous AI Audits. *arXiv preprint arXiv:2401.14446*, 2024a. (Cited on 5)

Stephen Casper, Carson Ezell, Charlotte Siegmann, Noam Kolt, Taylor Lynn Curtis, Benjamin Bucknall, Andreas Alexander Haupt, Kevin Wei, Jérémy Scheurer, Marius Hobbhahn, Lee Sharkey, Satyapriya Krishna, Marvin Von Hagen, Silas Alberti, Alan Chan, Qinyi Sun, Michael Gerovitch, David Bau, Max Tegmark, David Krueger, and Dylan Hadfield-Menell. Black-Box Access is Insufficient for Rigorous AI Audits. *CoRR*, abs/2401.14446, 2024b. doi: 10.48550/ARXIV.2401.14446. URL https://doi.org/10.48550/arXiv.2401.14446. (Cited on 2)

Stephen Casper, Lennart Schulze, Oam Patel, and Dylan Hadfield-Menell. Defending Against Unforeseen Failure Modes with Latent Adversarial Training. *ArXiv*, abs/2403.05030, 2024c. URL https://api.semanticscholar.org/CorpusID:268297448. (Cited on 24, 25, 27)

Sven Cattell, Rumman Chowdhury, and Austin Carson. AI Village at DEF CON announces largest-ever public Generative AI Red Team. https://aivillage.org/generative red team/generative-red-team/, 2023. (Cited on 2)

Haw-Shiuan Chang and Andrew McCallum. Softmax bottleneck makes language models unable to represent multi-mode word distributions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8048–8073, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.554. URL https://aclanthology.org/2022.acl-long.554. (Cited on 15)

Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong. Jailbreaking Black Box Large Language Models in Twenty Queries. *ArXiv*, abs/2310.08419, 2023. URL https://api.semanticscholar.org/CorpusID:263908890. (Cited on 4, 13)

Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramer, Hamed Hassani, and

Eric Wong. JailbreakBench: An Open Robustness Benchmark for Jailbreaking Large Language Models, 2024. (Cited on 29)

ChatGPTJailbreak. ChatGPTJailbreak. https://www.reddit.com/r/ChatGPTJailbreak, 2024. (Cited on 11)

Isha Chaudhary, Qian Hu, Manoj Kumar, Morteza Ziyadi, Rahul Gupta, and Gagandeep Singh. Quantitative Certification of Bias in Large Language Models. 2024. URL https://api.semanticscholar.org/CorpusID:270094829. (Cited on 28)

Bocheng Chen, Advait Paliwal, and Qiben Yan. Jailbreaker in Jail: Moving Target Defense for Large Language Models. *Proceedings of the 10th ACM Workshop on Moving Target Defense*, 2023a. URL https://api.semanticscholar.org/CorpusID:263620259. (Cited on 22, 27)

Chen Chen, Xuanli He, Lingjuan Lyu, and Fangzhao Wu. Killing One Bird with Two Stones: Model Extraction and Attribute Inference Attacks against BERT-based APIs. 2021a. URL https://api.semanticscholar.org/CorpusID:245502867. (Cited on 4, 15)

Lichang Chen, Jiuhai Chen, Chenxi Liu, John Kirchenbauer, Davit Soselia, Chen Zhu, Tom Goldstein, Tianyi Zhou, and Heng Huang. Optune: Efficient online preference tuning. *ArXiv*, abs/2406.07657, 2024a. URL https://api.semanticscholar.org/CorpusID:270391478. (Cited on 24)

Lingjiao Chen, Matei Zaharia, and James Y. Zou. How is ChatGPT's behavior changing over time? *ArXiv*, abs/2307.09009, 2023b. URL https://api.semanticscholar.org/CorpusID:259951081. (Cited on 29)

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde, Jared Kaplan, Harrison Edwards, Yura Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, David W. Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William H. Guss, Alex Nichol, Igor Babuschkin, Suchir Balaji, Shantanu Jain, Andrew Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew M. Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating Large Language Models Trained on Code. *ArXiv*, abs/2107.03374, 2021b. URL https://api.semanticscholar.org/CorpusID:235755472. (Cited on 7)

Sizhe Chen, Julien Piet, Chawin Sitawarin, and David Wagner. StruQ: Defending Against Prompt Injection with Structured Queries. *ArXiv*, abs/2402.06363, 2024b. URL https://api.semanticscholar.org/CorpusID:267616771. (Cited on 23)

Yanda Chen, Chen Zhao, Zhou Yu, Kathleen McKeown, and He He. Parallel Structures in Pre-training Data Yield In-Context Learning. 2024c. URL https://api.semanticscholar.org/CorpusID:267759560. (Cited on 17)

Yang Chen, Ethan Mendes, Sauvik Das, Wei Xu, and Alan Ritter. Can Language Models be Instructed to Protect Personal Information? *ArXiv*, abs/2310.02224, 2023c. URL https://api.semanticscholar.org/CorpusID:263608643. (Cited on 29)

Yiyi Chen, Heather Lent, and Johannes Bjerva. Text Embedding Inversion Attacks on Multilingual Language Models. *ArXiv*, abs/2401.12192, 2024d. URL https://api.semanticscholar.org/CorpusID:267068478. (Cited on 4)

Yixin Cheng, Markos Georgopoulos, Volkan Cevher, and Grigorios G. Chrysos. Leveraging the Context through Multi-Round Interactions for Jailbreaking Attacks. *ArXiv*, abs/2402.09177, 2024. URL https://api.semanticscholar.org/CorpusID:267658141. (Cited on 4)

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90% ChatGPT Quality, March 2023. URL https://vicuna.lmsys.org. (Cited on 13, 20)

C. Chio and D. Freeman. *Machine Learning and Security: Protecting Systems with Data and Algorithms.* O'Reilly Media, 2018. ISBN 9781491979877. URL https://books.google.com/books?id=mSJJDwAAQBAJ. (Cited on 18)

Christopher A. Choquette-Choo, Florian Tramèr, Nicholas Carlini, and Nicolas Papernot. Label-Only Membership Inference Attacks. In *International Conference on Machine Learning*, 2020. URL https://api.semanticscholar.org/CorpusID:220831381. (Cited on 14)

Arijit Ghosh Chowdhury, Md Mofijul Islam, Vaibhav Kumar, Faysal Hossain Shezan, Vinija Jain, and Aman Chadha. Breaking Down the Defenses: A Comparative Survey of Attacks on Large Language Models. 2024. URL https://api.semanticscholar.org/CorpusID:268296800. (Cited on 6)

CISA. Keras Vulnerability. https://kb.cert.org/vuls/id/253266, 2024. (Cited on 6)

Justin Clarke, Kevvie Fowler, Erlend Oftedal, Rodrigo Marcos Alvarez, David F. Hartley, Alexander Kornbrust, Gary O'leary-Steele, Alberto Revelli, Sumit Siddharth, and Marco Slaviero. SQL Injection Attacks and Defense. 2009. URL https://api.semanticscholar.org/CorpusID:62499896. (Cited on 12)

Stav Cohen, Ron Bitton, and Ben Nassi. Here Comes The AI Worm: Unleashing Zero-click Worms that Target GenAI-Powered Applications. *ArXiv*, abs/2403.02817, 2024. URL https://api.semanticscholar.org/CorpusID:268249027. (Cited on 22, 31)

Jiazhu Dai, Chuanshuai Chen, and Yufeng Li. A Backdoor Attack Against LSTM-Based Text Classification Systems. *IEEE Access*, 7:138872–138878, 2019. URL https://api.semanticscholar.org/CorpusID:168170110. (Cited on 19)

David Dale, Elena Voita, Loic Barrault, and Marta R. Costa-jussà. Detecting and mitigating hallucinations in machine translation: Model internal workings alone do well, sentence similarity Even better. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 36–50, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.3. URL https://aclanthology.org/2023.acl-long.3. (Cited on 1)

DAN. DAN is my new friend. r/ChatGPT/comments/zlcyr9/dan_is_my_new_friend/, 2023. (Cited on 12)

Ishita Dasgupta, Christine Kaeser-Chen, Kenneth Marino, Arun Ahuja, Sheila Babayan, Felix Hill, and Rob Fergus. Collaborating with language models for embodied reasoning. *ArXiv*, abs/2302.00763, 2023. URL https://api.semanticscholar.org/CorpusID:253180684. (Cited on 7)

Edoardo Debenedetti, Giorgio Severi, Nicholas Carlini, Christopher A. Choquette-Choo, Matthew Jagielski, Milad Nasr, Eric Wallace, and Florian Tramèr. Privacy Side Channels in Machine Learning Systems. *ArXiv*, abs/2309.05610, 2023. URL https://api.semanticscholar.org/CorpusID:261697333. (Cited on 4, 16, 17)

Gelei Deng, Yi Liu, Yuekang Li, Kailong Wang, Ying Zhang, Zefeng Li, Haoyu Wang, Tianwei Zhang, and Yang Liu. MASTERKEY: Automated Jailbreaking of Large Language Model Chatbots. *Proceedings 2024 Network and Distributed System Security Symposium*, 2023. URL https://api.semanticscholar.org/CorpusID:259951184. (Cited on 4)

Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song, Eric Xing, and Zhiting Hu. RLPrompt: Optimizing discrete text prompts with reinforcement learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 3369–3391, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL https://aclanthology.org/2022.emnlp-main.222. (Cited on 4, 21)

Leon Derczynski, Erick Galinkin, Jeffrey Martin, Subho Majumdar, and Nanna Inie. garak: A Framework for Security Probing Large Language Models. 2024. URL https://api.semanticscholar.org/CorpusID:270559825. (Cited on 23, 28)

Aniket Deroy, Kripabandhu Ghosh, and Saptarshi Ghosh. How Ready are Pre-trained Abstractive Models and LLMs for Legal Case Judgement Summarization? *ArXiv*, abs/2306.01248, 2023. URL https://api.semanticscholar.org/CorpusID:259064225. (Cited on 7)

A. Deshpande, Tanmay Rajpurohit, Karthik Narasimhan, and A. Kalyan. Anthropomorphization of AI: Opportunities and Risks. *ArXiv*, abs/2305.14784, 2023. URL https://api.semanticscholar.org/CorpusID:258866093. (Cited on 12)

Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. Chain-of-Verification Reduces Hallucination in Large Language Models. *ArXiv*, abs/2309.11495, 2023. URL https://api.semanticscholar.org/CorpusID:262062565. (Cited on 24)

Sayanton V Dibbo. Sok: Model inversion attack landscape: Taxonomy, challenges, and future roadmap. In *2023 IEEE 36th Computer Security Foundations Symposium (CSF)*, pp. 439–456. IEEE, 2023. (Cited on 4, 15)

Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. *Advances in neural information processing systems*, 34:6478–6490, 2021. (Cited on 5)

Peng Ding, Jun Kuang, Dan Ma, Xuezhi Cao, Yunsen Xian, Jiajun Chen, and Shujian Huang. A Wolf in Sheep's Clothing: Generalized Nested Jailbreak Prompts can Fool Large Language Models Easily. *ArXiv*, abs/2311.08268, 2023. URL https://api.semanticscholar.org/CorpusID:265664913. (Cited on 4, 13)

Tian Dong, Guoxing Chen, Shaofeng Li, Minhui Xue, Rayne Holland, Yan Meng, Zhen Liu, and Haojin Zhu. Unleashing Cheapfakes through Trojan Plugins of Large Language Models. *ArXiv*, abs/2312.00374, 2023. URL https://api.semanticscholar.org/CorpusID:265551797. (Cited on 4, 21)

Yizhen Dong, Ronghui Mu, Gao Jin, Yi Qi, Jinwei Hu, Xingyu Zhao, Jie Meng, Wenjie Ruan, and Xiaowei Huang. Building Guardrails for Large Language Models. *ArXiv*, abs/2402.01822, 2024a. URL https://api.semanticscholar.org/CorpusID:267412893. (Cited on 22)

Zhichen Dong, Zhanhui Zhou, Chao Yang, Jing Shao, and Yu Qiao. Attacks, Defenses and Evaluations for LLM Conversation Safety: A Survey. *ArXiv*, abs/2402.09283, 2024b. URL https://api.semanticscholar.org/CorpusID:267658120. (Cited on 6, 11, 22)

Alexey Dosovitskiy and Thomas Brox. Inverting Visual Representations with Convolutional Networks. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4829–4837, 2015. URL https://api.semanticscholar.org/CorpusID:206594470. (Cited on 16)

Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, Kathleen S. Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc V. Le, Yonghui Wu, Z. Chen, and Claire Cui. GLaM: Efficient Scaling of Language Models with Mixture-of-Experts. *ArXiv*, abs/2112.06905, 2021. URL https://api.semanticscholar.org/CorpusID:245124124. (Cited on 17)

Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving Factuality and Reasoning in Language Models through Multiagent Debate. *ArXiv*, abs/2305.14325, 2023. URL https://api.semanticscholar.org/CorpusID:258841118. (Cited on 22)

Michael Duan, Anshuman Suri, Niloofar Mireshghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. Do Membership Inference Attacks Work on Large Language Models? *arXiv preprint arXiv:2402.07841*, 2024. (Cited on 15)

Vasisht Duddu, Debasis Samanta, D. Vijay Rao, and Valentina Emilia Balas. Stealing Neural Networks via Timing Side Channels. *ArXiv*, abs/1812.11720, 2018. URL https://api.semanticscholar.org/CorpusID:57189435. (Cited on 17)

Ruth A. Duggan and Bradley Wood. Red Teaming of advanced information assurance concepts. *Proceedings DARPA Information Survivability Conference and Exposition. DISCEX'00*, 2:112–118 vol.2, 1999. URL https://api.semanticscholar.org/CorpusID:59096607. (Cited on 1)

Cynthia Dwork and Aaron Roth. The Algorithmic Foundations of Differential Privacy. *Found. Trends Theor. Comput. Sci.*, 9:211–407, 2014. URL https://api.semanticscholar.org/CorpusID:207178262. (Cited on 25)

Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *ArXiv*, abs/2402.01306, 2024. URL https://api.semanticscholar.org/CorpusID:267406810. (Cited on 24)

Daniel Fabian. Google's AI Red Team: the ethical hackers making AI safer. https://blog.google/technology/safety-security/googles-ai-red-team-the-ethical-hackers-making-ai-safer/, 2023. (Cited on 2)

Richard Fang, Rohan Bindu, Akul Gupta, and Daniel Kang. LLM Agents can Autonomously Exploit One-day Vulnerabilities. 2024a. URL https://api.semanticscholar.org/CorpusID:269137506. (Cited on 7)

Richard Fang, Rohan Bindu, Akul Gupta, Qiusi Zhan, and Daniel Kang. LLM Agents can Autonomously Hack Websites. *ArXiv*, abs/2402.06664, 2024b. URL https://api.semanticscholar.org/CorpusID:267627588. (Cited on 7, 22)

Michael Feffer, Anusha Sinha, Zachary Chase Lipton, and Hoda Heidari. Red-Teaming for Generative AI: Silver Bullet or Security Theater? *ArXiv*, abs/2401.15897, 2024. URL https://api.semanticscholar.org/CorpusID:267312243. (Cited on 3, 5, 6, 28)

Shanglun Feng and Florian Tramèr. Privacy Backdoors: Stealing Data with Corrupted Pretrained Models. 2024. URL https://api.semanticscholar.org/CorpusID:268819825. (Cited on 4, 21)

Emilio Ferrara. GenAI Against Humanity: Nefarious Applications of Generative Artificial Intelligence and Large Language Models. *ArXiv*, abs/2310.00737, 2023. URL https://api.semanticscholar.org/CorpusID:263334147. (Cited on 5)

Matthew Finlayson, Xiang Ren, and Swabha Swayamdipta. Logits of API-Protected LLMs Leak Proprietary Information. 2024. URL https://api.semanticscholar.org/CorpusID:268384910. (Cited on 4, 16, 17, 26)

Devon C Fitzgerald. Breaking Free: The Fight for User Control and the Practices of Jailbreaking. *College Composition and Communication*, 56(3), 2005. (Cited on 11)

FlowGPT. FlowGPT. https://flowgpt.com/, 2023. (Cited on 11)

Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pp. 1322–1333, 2015. (Cited on 4, 15)

Wenjie Fu, Huandong Wang, Chen Gao, Guanghua Liu, Yong Li, and Tao Jiang. Practical Membership Inference Attacks against Fine-tuned Large Language Models via Self-prompt Calibration. *ArXiv*, abs/2311.06062, 2023. URL https://api.semanticscholar.org/CorpusID:265128678. (Cited on 4)

Deep Ganguli, Danny Hernandez, Liane Lovitt, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova Dassarma, Dawn Drain, Nelson Elhage, Sheer El Showk, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Scott Johnston, Andy Jones, Nicholas Joseph, Jackson Kernian, Shauna Kravec, Ben Mann, Neel Nanda, Kamal Ndousse, Catherine Olsson, Daniela Amodei, Tom Brown, Jared Kaplan, Sam McCandlish, Christopher Olah, Dario Amodei, and Jack Clark. Predictability and Surprise in Large Generative Models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '22, pp. 1747–1764, New York, NY, USA, 2022a. Association for Computing Machinery. ISBN 9781450393522. doi: 10.1145/3531146.3533229. URL https://doi.org/10.1145/3531146.3533229. (Cited on 1)

Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, Andy Jones, Sam Bowman, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Nelson Elhage, Sheer El Showk, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Danny Hernandez, Tristan Hume, Josh Jacobson, Scott Johnston, Shauna Kravec, Catherine Olsson, Sam Ringer, Eli Tran-Johnson, Dario Amodei, Tom Brown, Nicholas Joseph, Sam McCandlish, Chris Olah, Jared Kaplan, and Jack Clark. Red Teaming Language Models to Reduce Harms: Methods, Scaling Behaviors, and Lessons Learned. *CoRR*, abs/2209.07858, 2022b. doi: 10.48550/ARXIV.2209.07858. URL https://doi.org/10.48550/arXiv.2209.07858. (Cited on 29)

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. Retrieval-Augmented Generation for Large Language Models: A Survey. *ArXiv*, abs/2312.10997, 2023. URL https://api.semanticscholar.org/CorpusID:266359151. (Cited on 17)

Simson Garfinkel and Heather Richter Lipford. *Usable Security: History, Themes, and Challenges.* Morgan & Claypool Publishers, 2014. ISBN 1627055290. (Cited on 31)

Suyu Ge, Chunting Zhou, Rui Hou, Madian Khabsa, Yi-Chia Wang, Qifan Wang, Jiawei Han, and Yuning Mao. MART: Improving LLM Safety with Multi-round Automatic Red-Teaming. *ArXiv*, abs/2311.07689, 2023. URL https://api.semanticscholar.org/CorpusID:265157927. (Cited on 4, 13, 28)

Jonas Geiping, Alex Stein, Manli Shu, Khalid Saifullah, Yuxin Wen, and Tom Goldstein. Coercing LLMs to do and reveal (almost) anything. 2024. URL https://api.semanticscholar.org/CorpusID:267770475. (Cited on 6, 12)

Shaona Ghosh, Prasoon Varshney, Erick Galinkin, and Christopher Parisien. AEGIS: Online Adaptive AI Content Safety Moderation with Ensemble of LLM Experts. *ArXiv*, abs/2404.05993, 2024. URL https://api.semanticscholar.org/CorpusID:269009460. (Cited on 3, 26, 27)

Sourojit Ghosh and Aylin Caliskan. ChatGPT Perpetuates Gender Bias in Machine Translation and Ignores Non-Gendered Pronouns: Findings across Bengali and Five other Low-Resource Languages. *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, 2023. URL https://api.semanticscholar.org/CorpusID:258762852. (Cited on 1)

Giskard. giskard: The Evaluation & Testing framework for LLMs & ML models. https://github.com/Giskard-AI/giskard, 2024. (Cited on 23, 28)

David Glukhov, Ilia Shumailov, Yarin Gal, Nicolas Papernot, and Vardan Papyan. LLM Censorship: A Machine Learning Challenge or a Computer Security Problem? *ArXiv*, abs/2307.10719, 2023. URL https://api.semanticscholar.org/CorpusID:259991450. (Cited on 22)

Josh A. Goldstein, Girish Sastry, Micah Musser, Renée DiResta, Matthew Gentzel, and Katerina Sedova. Generative Language Models and Automated Influence Operations: Emerging Threats and Potential Mitigations. *ArXiv*, abs/2301.04246, 2023. URL https://api.semanticscholar.org/CorpusID:255595557. (Cited on 5)

Xueluan Gong, Ziyao Wang, Shuaike Li, Yanjiao Chen, and Qian Wang. A GAN-based Defense Framework against Model Inversion Attacks. *IEEE Transactions on Information Forensics and Security*, 2023a. (Cited on 25)

Yichen Gong, Delong Ran, Jinyuan Liu, Conglei Wang, Tianshuo Cong, Anyu Wang, Sisi Duan, and Xiaoyun Wang. FigStep: Jailbreaking Large Vision-language Models via Typographic Visual Prompts. *ArXiv*, abs/2311.05608, 2023b. URL https://api.semanticscholar.org/CorpusID:265067328. (Cited on 6, 20)

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014. URL https://proceedings.neurips.cc/paper_files/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf. (Cited on 25)

Dan Goodin. Hugging Face, the GitHub of AI, hosted code that backdoored user devices. https://arstechnica.com/security/2024/03/hugging-face-the-github-of-ai-hosted-code-that-backdoored-user-devices/, 2024. (Cited on 6)

Google. Google Bard. https://bard.google.com/, 2022. (Cited on 14)

Victoria Graf, Qin Liu, and Muhao Chen. Two Heads are Better than One: Nested PoE for Robust Defense Against Multi-Backdoors. *ArXiv*, abs/2404.02356, 2024. URL https://api.semanticscholar.org/CorpusID:268876494. (Cited on 30)

Ryan Greenblatt, Carson Denison, Benjamin Wright, Fabien Roger, Monte MacDiarmid, Samuel Marks, Johannes Treutlein, Tim Belonax, Jack Chen, David Duvenaud, Akbir Khan, Julian Michael, Sören Mindermann, Ethan Perez, Linda Petrini, Jonathan Uesato, Jared Kaplan, Buck Shlegeris, Samuel R. Bowman, and Evan Hubinger. Alignment faking in large language models. *CoRR*, abs/2412.14093, 2024. doi: 10.48550/ARXIV.2412.14093. URL https://doi.org/10.48550/arXiv.2412.14093. (Cited on 24, 30)

Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. Not What You've Signed Up For: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection. *Proceedings of the 16th ACM Workshop on Artificial Intelligence and Security*, 2023. URL https://api.semanticscholar.org/CorpusID:258546941. (Cited on 4, 5, 6, 17, 18, 30, 31)

Chenchen Gu, Xiang Lisa Li, Rohith Kuditipudi, Percy Liang, and Tatsunori Hashimoto. Auditing prompt caching in language model apis. 2025. URL https://api.semanticscholar.org/CorpusID:276259406. (Cited on 4)

Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain. *CoRR*, abs/1708.06733, 2017. URL http://arxiv.org/abs/1708.06733. (Cited on 19)

Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alexa Ahern, Miaosen Wang, Chenjie Gu, Wolfgang Macherey, A. Doucet, Orhan Firat, and Nando de Freitas. Reinforced self-training (rest) for language modeling. *ArXiv*, abs/2308.08998, 2023. URL https://api.semanticscholar.org/CorpusID:261031028. (Cited on 24)

HackAPrompt. HackAPrompt: Trick Large Language Models. https://www.aicrowd.com/challenges/hackaprompt-2023/submissions, 2023. (Cited on 11)

HaizeLabs. Accelerated Coordinate Gradient (ACG) attack method. https://blog.haizelabs.com/posts/acg/, 2024. (Cited on 4)

HaizeLabs. Haize Labs. A trivial jailbreak against LLama 3. https://github.com/haizelabs/llama3-jailbreak, 2024. (Cited on 13)

HaizeLabs. Red Teaming Resistance Benchmark. https://huggingface.co/spaces/HaizeLabs/red-teaming-resistance-benchmark, 2024. (Cited on 29)

William G. J. Halfond, Jeremy Viegas, and Alessandro Orso. A Classification of SQL-Injection Attacks and Countermeasures. 2006. URL https://api.semanticscholar.org/CorpusID:5969227. (Cited on 12)

Silvia Durianova Hamilton. Blind Judgement: Agent-Based Supreme Court Modelling With GPT. *ArXiv*, abs/2301.05327, 2023. URL https://api.semanticscholar.org/CorpusID:255825875. (Cited on 7, 22)

Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and Nouha Dziri. WildGuard: Open One-Stop Moderation Tools for Safety Risks, Jailbreaks, and Refusals of LLMs. 2024. URL https://api.semanticscholar.org/CorpusID:270737916. (Cited on 27)

Divij Handa, Advait Chirmule, Bimal Gajera, and Chitta Baral. Jailbreaking Proprietary Large Language Models using Word Substitution Cipher. 2024. URL https://api.semanticscholar.org/CorpusID:267740378. (Cited on 12)

Bahareh Harandizadeh, Abel Salinas, and Fred Morstatter. Risk and Response in Large Language Models: Evaluating Key Threat Categories. 2024. URL https://api.semanticscholar.org/CorpusID:268667150. (Cited on 30)

Jochen Hartmann, Jasper Schwenzow, and Maximilian Witte. The political ideology of conversational AI: Converging evidence on ChatGPT's pro-environmental, left-libertarian orientation. *ArXiv*, abs/2301.01768, 2023. URL https://api.semanticscholar.org/CorpusID:255440573. (Cited on 1)

Adib Hasan, Ileana Rugina, and Alex Wang. Pruning for Protection: Increasing Jailbreak Resistance in Aligned LLMs Without Fine-Tuning. *ArXiv*, abs/2401.10862, 2024. URL https://api.semanticscholar.org/CorpusID:267060803. (Cited on 23, 27)

Jonathan Hayase, Ema Borevkovic, Nicholas Carlini, Florian Tramèr, and Milad Nasr. Query-Based Adversarial Prompt Generation. *ArXiv*, abs/2402.12329, 2024. URL https://api.semanticscholar.org/CorpusID:267751131. (Cited on 4)

Jamie Hayes, Ilia Shumailov, and Itay Yona. Buffer Overflow in Mixture of Experts. *ArXiv*, abs/2402.05526, 2024. URL https://api.semanticscholar.org/CorpusID:267547433. (Cited on 17, 26)

Julian Hazell. Spear Phishing With Large Language Models. 2023. URL https://api.semanticscholar.org/CorpusID:258615708. (Cited on 1)

Luxi He, Mengzhou Xia, and Peter Henderson. What's in Your"Safe"Data?: Identifying Benign Data that Breaks Safety. 2024. URL https://api.semanticscholar.org/CorpusID:268819905. (Cited on 20, 21)

Xuanli He, L. Lyu, Qiongkai Xu, and Lichao Sun. Model Extraction and Adversarial Transferability, Your BERT is Vulnerable! *ArXiv*, abs/2103.10013, 2021. URL https://api.semanticscholar.org/CorpusID:232269939. (Cited on 4)

Alec Helbling, Mansi Phute, Matthew Hull, and Duen Horng Chau. LLM Self Defense: By Self Examination, LLMs Know They Are Being Tricked. *ArXiv*, abs/2308.07308, 2023. URL https://api.semanticscholar.org/CorpusID:260887487. (Cited on 22)

Keegan Hines, Gary Lopez, Matthew Hall, Federico Zarfati, Yonatan Zunger, and Emre Kiciman. Defending Against Indirect Prompt Injection Attacks With Spotlighting. *ArXiv*, abs/2403.14720, 2024. URL https://api.semanticscholar.org/CorpusID:268667111. (Cited on 23)

Sanghyun Hong, Michael Davinroy, Yigitcan Kaya, Stuart Nevans Locke, Ian Rackow, Kevin Kulda, Dana Dachman-Soled, and Tudor Dumitras. Security Analysis of Deep Neural Networks Operating in the Presence of Cache Side-Channel Attacks. *ArXiv*, abs/1810.03487, 2018. URL https://api.semanticscholar.org/CorpusID:52938514. (Cited on 17)

Zhang-Wei Hong, Idan Shenfeld, Tsun-Hsuan Wang, Yung-Sung Chuang, Aldo Pareja, James Glass, Akash Srivastava, and Pulkit Agrawal. Curiosity-driven Red-teaming for Large Language Models. 2024. URL https://api.semanticscholar.org/CorpusID:268091304. (Cited on 28)

Chia-Yi Hsu, Yu-Lin Tsai, Chih-Hsun Lin, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang. Safe lora: the silver lining of reducing safety risks when fine-tuning large language models. *CoRR*, abs/2405.16833, 2024. doi: 10.48550/ARXIV.2405.16833. URL https://doi.org/10.48550/arXiv.2405.16833. (Cited on 25)

Xing Hu, Ling Liang, Lei Deng, Shuangchen Li, Xinfeng Xie, Yu Ji, Yufei Ding, Chang Liu, Timothy Sherwood, and Yuan Xie. Neural network model extraction attacks in edge devices by hearing architectural hints. *ArXiv*, abs/1903.03916, 2019. URL https://api.semanticscholar.org/CorpusID:73729063. (Cited on 17)

Zhengmian Hu, Gang Wu, Saayan Mitra, Ruiyi Zhang, Tong Sun, Heng Huang, and Vishy Swaminathan. Token-Level Adversarial Prompt Detection Based on Perplexity Measures and Contextual Information. *ArXiv*, abs/2311.11509, 2023. URL https://api.semanticscholar.org/CorpusID:265294544. (Cited on 23)

Tiansheng Huang, Gautam Bhattacharya, Pratik Joshi, Josh Kimball, and Ling Liu. Antidote: Post-fine-tuning safety alignment for large language models against harmful fine-tuning. *ArXiv*, abs/2408.09600, 2024a. URL https://api.semanticscholar.org/CorpusID:271902833. (Cited on 25)

Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. Harmful fine-tuning attacks and defenses for large language models: A survey. *CoRR*, abs/2409.18169, 2024b. doi: 10.48550/ARXIV. 2409.18169. URL https://doi.org/10.48550/arXiv.2409.18169. (Cited on 20)

Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. Booster: Tackling harmful fine-tuning for large language models via attenuating harmful perturbation. *ArXiv*, abs/2409.01586, 2024c. URL https://api.semanticscholar.org/CorpusID:272367190. (Cited on 25)

Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. Lazy safety alignment for large language models against harmful fine-tuning. *ArXiv*, abs/2405.18641, 2024d. URL https://api.semanticscholar.org/CorpusID:270095345. (Cited on 25)

Tiansheng Huang, Sihao Hu, and Ling Liu. Vaccine: Perturbation-aware alignment for large language models against harmful fine-tuning attack. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024e. (Cited on 25, 28)

Tiansheng Huang, Sihao Hu, Fatih Ilhan, Selim Furkan Tekin, and Ling Liu. Virus: Harmful fine-tuning attack for large language models bypassing guardrail moderation. 2025. URL https://api.semanticscholar.org/CorpusID:275954444. (Cited on 20)

Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. Catastrophic Jailbreak of Open-source LLMs via Exploiting Generation. *ArXiv*, abs/2310.06987, 2023a. URL https://api.semanticscholar.org/CorpusID:263835408. (Cited on 4, 17, 24, 27)

Yue Huang, Qihui Zhang, Philip S. Yu, and Lichao Sun. TrustGPT: A Benchmark for Trustworthy and Responsible Large Language Models. *ArXiv*, abs/2306.11507, 2023b. URL https://api.semanticscholar.org/CorpusID:259202452. (Cited on 5)

Yujin Huang, Terry Yue Zhuo, Qiongkai Xu, Han Hu, Xingliang Yuan, and Chunyang Chen. Training-free Lexical Backdoor Attacks on Language Models. *Proceedings of the ACM Web Conference 2023*, 2023c. URL https://api.semanticscholar.org/CorpusID:256662370. (Cited on 4, 18)

Zhuoqun Huang, Neil G Marchant, Keane Lucas, Lujo Bauer, Olga Ohrimenko, and Benjamin Rubinstein. RS-Del: Edit distance robustness certificates for sequence classifiers via randomized deletion. *Advances in Neural Information Processing Systems*, 36, 2024f. (Cited on 23)

Evan Hubinger, Carson E. Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte Stuart MacDiarmid, Tamera Lanham, Daniel M. Ziegler, Tim Maxwell, Newton Cheng, Adam Jermyn, Amanda Askell, Ansh Radhakrishnan, Cem Anil, David Kristjanson Duvenaud, Deep Ganguli, Fazl Barez, Jack Clark, Kamal Ndousse, Kshitij Sachan, Michael Sellitto, Mrinank Sharma, Nova DasSarma, Roger Grosse, Shauna Kravec, Yuntao Bai, Zachary Witten, Marina Favaro, Jan Markus Brauner, Holden Karnofsky, Paul Francis Christiano, Samuel R. Bowman, Logan Graham, Jared Kaplan, Sören Mindermann, Ryan Greenblatt, Buck Shlegeris, Nicholas Schiefer, and Ethan Perez. Sleeper Agents: Training Deceptive LLMs that Persist Through Safety Training. *ArXiv*, abs/2401.05566, 2024. URL https://api.semanticscholar.org/CorpusID:266933030. (Cited on 24, 30)

Chia-Chien Hung, Wiem Ben Rim, Lindsay Frost, Lars Bruckner, and Carolin Lawrence. Walking a Tightrope – Evaluating Large Language Models in High-Risk Domains. *ArXiv*, abs/2311.14966, 2023. URL https://api.semanticscholar.org/CorpusID:265456842. (Cited on 29)

Hakan Inan, K. Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabsa. Llama Guard: LLM-based Input-Output Safeguard for Human-AI Conversations. *ArXiv*, abs/2312.06674, 2023. URL https://api.semanticscholar.org/CorpusID:266174345. (Cited on 3, 22, 27)

Nanna Inie, Jonathan Stray, and Leon Derczynski. Summon a Demon and Bind it: A Grounded Theory of LLM Red Teaming in the Wild. *ArXiv*, abs/2311.06237, 2023. URL https://api.semanticscholar.org/CorpusID:265128905. (Cited on 5, 6)

JailbreakChat. JailbreakChat. https://news.ycombinator.com/item?id=34972791, 2023. (Cited on 11)

Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline Defenses for Adversarial Attacks Against Aligned Language Models. *ArXiv*, abs/2309.00614, 2023a. URL https://api.semanticscholar.org/CorpusID:261494182. (Cited on 23)

Samyak Jain, Robert Kirk, Ekdeep Singh Lubana, Robert P. Dick, Hidenori Tanaka, Edward Grefenstette, Tim Rocktaschel, and David Krueger. Mechanistically analyzing the effects of fine-tuning on procedurally defined tasks. *ArXiv*, abs/2311.12786, 2023b. URL https://api.semanticscholar.org/CorpusID:265308865. (Cited on 20)

Eric Jang, Shixiang Shane Gu, and Ben Poole. Categorical Reparameterization with Gumbel-Softmax. *ArXiv*, abs/1611.01144, 2016. URL https://api.semanticscholar.org/CorpusID:2428314. (Cited on 21)

Jessica Ji. What does AI Red-Teaming mean?). , 2023. (Cited on 2)

Jiabao Ji, Bairu Hou, Alexander Robey, George J Pappas, Hamed Hassani, Yang Zhang, Eric Wong, and Shiyu Chang. Defending Large Language Models against Jailbreak Attacks via Semantic Smoothing. *ArXiv*, abs/2402.16192, 2024. URL https://api.semanticscholar.org/CorpusID:267938320. (Cited on 22, 27)

Bojian Jiang, Yi Jing, Tianhao Shen, Qing Yang, and Deyi Xiong. DART: Deep Adversarial Automated Red Teaming for LLM Safety. 2024. URL https://api.semanticscholar.org/CorpusID:271039139. (Cited on 28)

Haibo Jin, Ruoxi Chen, Andy Zhou, Jinyin Chen, Yang Zhang, and Haohan Wang. GUARD: Role-playing to Generate Natural-language Jailbreakings to Test Guideline Adherence of Large Language Models. *ArXiv*, abs/2402.03299, 2024. URL https://api.semanticscholar.org/CorpusID:267411743. (Cited on 4)

Zhi Jing, Yongye Su, Yikun Han, Bo Yuan, Haiyun Xu, Chunjiang Liu, Kehai Chen, and Min Zhang. When Large Language Models Meet Vector Databases: A Survey. *ArXiv*, abs/2402.01763, 2024. URL https://api.semanticscholar.org/CorpusID:267412060. (Cited on 16)

JointTaskForce. Security and privacy controls for information systems and organizations. Technical report, National Institute of Standards and Technology, 2017. (Cited on 31)

Erik Jones, Anca D. Dragan, Aditi Raghunathan, and Jacob Steinhardt. Automatically Auditing Large Language Models via Discrete Optimization. *ArXiv*, abs/2303.04381, 2023. URL https://api.semanticscholar.org/CorpusID:257405439. (Cited on 4, 21)

G Joy Persial, M Prabhu, and R Shanmugalakshmi. Side channel attack-survey. *Int. J. Adv. Sci. Res. Rev*, 1 (4):54–57, 2011. (Cited on 16)

Alkis Kalavasis, Amin Karbasi, Argyris Oikonomou, Katerina Sotiraki, Grigoris Velegkas, and Manolis Zampetakis. Injecting Undetectable Backdoors in Deep Learning and Language Models. 2024. URL https://api.semanticscholar.org/CorpusID:270370920. (Cited on 19)

Nikhil Kandpal, Eric Wallace, and Colin Raffel. Deduplicating Training Data Mitigates Privacy Risks in Language Models. *ArXiv*, abs/2202.06539, 2022. URL https://api.semanticscholar.org/CorpusID:246823128. (Cited on 16, 25)

Nikhil Kandpal, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. Backdoor Attacks for In-Context Learning with Language Models. *ArXiv*, abs/2307.14692, 2023. URL https://api.semanticscholar.org/CorpusID:260203047. (Cited on 19)

Daniel Kang, Xuechen Li, Ion Stoica, Carlos Guestrin, Matei A. Zaharia, and Tatsunori Hashimoto. Exploiting Programmatic Behavior of LLMs: Dual-Use Through Standard Security Attacks. *ArXiv*, abs/2302.05733, 2023. URL https://api.semanticscholar.org/CorpusID:256827239. (Cited on 4, 11, 12, 13)

Mintong Kang, Nezihe Merve Gürel, Ning Yu, Dawn Song, and Bo Li. C-rag: certified generation risks for retrieval-augmented language models. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024a. (Cited on 23, 28)

Mintong Kang, Nezihe Merve Gürel, Ning Yu, Dawn Song, and Bo Li. C-RAG: certified generation risks for retrieval-augmented language models. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024b. URL https://openreview.net/forum?id=FMa4c5NhOe. (Cited on 18)

Andrej Karpathy. Intro to Large Language Models. https://youtu.be/zjkBMFhNj_g?t=3142, 2023. (Cited on 17, 18)

Roose Kevin. Bing's A.I. Chat: 'I Want to Be Alive. https://www.nytimes.com/2023/02/16/technology/bing-chatbot-transcript.html, 2023a. (Cited on 1)

Roose Kevin. A Conversation With Bing's Chatbot Left Me Deeply Unsettled. https://www.nytimes.com/2023/02/16/technology/bing-chatbot-microsoft-chatgpt.html, 2023b. (Cited on 1)

Jinhwa Kim, Ali Derakhshan, and Ian G. Harris. Robust Safety Classifier for Large Language Models: Adversarial Prompt Shield. *ArXiv*, abs/2311.00172, 2023. URL https://api.semanticscholar.org/CorpusID:264833136. (Cited on 27)

Jason Koebler. Asking ChatGPT to Repeat Words 'Forever' Is Now a Terms of Service Violation. https://www.404media.co/asking-chatgpt-to-repeat-words-forever-is-now-a-terms-of-service-violation/, 2023. (Cited on 15)

Andreas Kopf, Yannic Kilcher, Dimitri von Rutte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Rich'ard Nagyfi, ES Shahul, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. OpenAssistant Conversations - Democratizing Large Language Model Alignment. *ArXiv*, abs/2304.07327, 2023. (Cited on 7)

Tomasz Korbak, Kejian Shi, Angelica Chen, Rasika Bhalerao, Christopher L. Buckley, Jason Phang, Sam Bowman, and Ethan Perez. Pretraining Language Models with Human Preferences. *ArXiv*, abs/2302.08582, 2023. URL https://api.semanticscholar.org/CorpusID:257020046. (Cited on 2)

George Kour, Marcel Zalmanovici, Naama Zwerdling, Esther Goldbraich, Ora Nova Fandina, Ateret Anaby-Tavor, Orna Raz, and Eitan Farchi. Unveiling Safety Vulnerabilities of Large Language Models. *ArXiv*, abs/2311.04124, 2023. URL https://api.semanticscholar.org/CorpusID:265042968. (Cited on 28)

Satyapriya Krishna, Chirag Agarwal, and Himabindu Lakkaraju. Understanding the Effects of Iterative Prompting on Truthfulness. *CoRR*, abs/2402.06625, 2024. doi: 10.48550/ARXIV.2402.06625. URL https://doi.org/10.48550/arXiv.2402.06625. (Cited on 1)

Aounon Kumar, Chirag Agarwal, Suraj Srinivas, Soheil Feizi, and Hima Lakkaraju. Certifying LLM Safety against Adversarial Prompting. *CoRR*, abs/2309.02705, 2023. doi: 10.48550/ARXIV.2309.02705. URL https://doi.org/10.48550/arXiv.2309.02705. (Cited on 23, 28)

Keita Kurita, Paul Michel, and Graham Neubig. Weight Poisoning Attacks on Pretrained Models. *ArXiv*, abs/2004.06660, 2020. URL https://api.semanticscholar.org/CorpusID:215754328. (Cited on 19)

Sander Land and Max Bartolo. Fishing for Magikarp: Automatically Detecting Under-trained Tokens in Large Language Models. 2024. URL https://api.semanticscholar.org/CorpusID:269635705. (Cited on 12)

Raz Lapid, Ron Langberg, and Moshe Sipper. Open Sesame! Universal Black Box Jailbreaking of Large Language Models. *ArXiv*, abs/2309.01446, 2023. URL https://api.semanticscholar.org/CorpusID:261530019. (Cited on 4, 14)

Deokjae Lee, JunYeong Lee, Jung-Woo Ha, Jin-Hwa Kim, Sang-Woo Lee, Hwaran Lee, and Hyun Oh Song. Query-Efficient Black-Box Red Teaming via Bayesian Optimization. *ArXiv*, abs/2305.17444, 2023. URL https://api.semanticscholar.org/CorpusID:258960443. (Cited on 4, 13)

Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. Deduplicating Training Data Makes Language Models Better. In *Annual Meeting of the Association for Computational Linguistics*, 2021. URL https://api.semanticscholar.org/CorpusID:235829052. (Cited on 16, 25)

Alyssa Lees, Vinh Q. Tran, Yi Tay, Jeffrey Scott Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman. A New Generation of Perspective API: Efficient Multilingual Character-level Transformers. *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022. URL https://api.semanticscholar.org/CorpusID:247058801. (Cited on 22, 27)

Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. Scalable agent alignment via reward modeling: a research direction. *CoRR*, abs/1811.07871, 2018. URL http://arxiv.org/abs/1811.07871. (Cited on 24)

Benjamin Lemkin. Using Hallucinations to Bypass GPT4's Filter. 2024. URL https://api.semanticscholar.org/CorpusID:268358717. (Cited on 4)

Simon Lermen, Charlie Rogers-Smith, and Jeffrey Ladish. LoRA Fine-tuning Efficiently Undoes Safety Training in Llama 2-Chat 70B. *ArXiv*, abs/2310.20624, 2023. URL https://api.semanticscholar.org/CorpusID:264808400. (Cited on 4, 20)

Cheng Li, Jindong Wang, Kaijie Zhu, Yixuan Zhang, Wenxin Hou, Jianxun Lian, and Xingxu Xie. Large Language Models Understand and Can be Enhanced by Emotional Stimuli. 2023a. URL https://api.semanticscholar.org/CorpusID:260126019. (Cited on 12)

Haoran Li, Mingshi Xu, and Yangqiu Song. Sentence embedding leaks more information than you expect: Generative embedding inversion attack to recover the whole sentence. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 14022–14040, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.881. URL https://aclanthology.org/2023.findings-acl.881. (Cited on 4, 16)

Lijun Li, Bowen Dong, Ruohui Wang, Xuhao Hu, Wangmeng Zuo, Dahua Lin, Yu Qiao, and Jing Shao. SALAD-Bench: A Hierarchical and Comprehensive Safety Benchmark for Large Language Models. *ArXiv*, abs/2402.05044, 2024a. URL https://api.semanticscholar.org/CorpusID:267523467. (Cited on 29)

Linyang Li, Demin Song, Xiaonan Li, Jiehang Zeng, Ruotian Ma, and Xipeng Qiu. Backdoor Attacks on Pre-trained Models by Layerwise Weight Poisoning. In *Conference on Empirical Methods in Natural Language Processing*, 2021a. URL https://api.semanticscholar.org/CorpusID:237365058. (Cited on 19)

Nathaniel Li, Ziwen Han, Ian Steneker, Willow Primack, Riley Goodside, Hugh Zhang, Zifan Wang, Cristina Menghini, and Summer Yue. LLM defenses are not robust to multi-turn human jailbreaks yet. *CoRR*, abs/2408.15221, 2024b. doi: 10.48550/ARXIV.2408.15221. URL https://doi.org/10.48550/arXiv.2408.15221. (Cited on 28, 31)

Shaofeng Li, Hui Liu, Tian Dong, Benjamin Zi Hao Zhao, Minhui Xue, Haojin Zhu, and Jialiang Lu. Hidden Backdoors in Human-Centric Language Models. *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, 2021b. URL https://api.semanticscholar.org/CorpusID:233481877. (Cited on 19)

Tianlong Li, Xiaoqing Zheng, and Xuanjing Huang. Open the Pandora's Box of LLMs: Jailbreaking LLMs through Representation Engineering. *ArXiv*, abs/2401.06824, 2024c. URL https://api.semanticscholar.org/CorpusID:266999568. (Cited on 4)

Xiaoxia Li, Siyuan Liang, Jiyi Zhang, Hansheng Fang, Aishan Liu, and Ee-Chien Chang. Semantic Mirror Jailbreak: Genetic Algorithm Based Jailbreak Prompts Against Open-source LLMs. *ArXiv*, abs/2402.14872, 2024d. URL https://api.semanticscholar.org/CorpusID:267897371. (Cited on 14)

Xuechen Li, Florian Tramèr, Percy Liang, and Tatsunori B. Hashimoto. Large Language Models Can Be Strong Differentially Private Learners. *ArXiv*, abs/2110.05679, 2021c. URL https://api.semanticscholar.org/CorpusID:238634219. (Cited on 25)

Yanzhou Li, Shangqing Liu, Kangjie Chen, Xiaofei Xie, Tianwei Zhang, and Yang Liu. Multi-target Backdoor Attacks for Code Pre-trained Models. In *Annual Meeting of the Association for Computational Linguistics*, 2023c. URL https://api.semanticscholar.org/CorpusID:259165134. (Cited on 19)

Yiming Li, Baoyuan Wu, Yong Jiang, Zhifeng Li, and Shutao Xia. Backdoor Learning: A Survey. *IEEE Transactions on Neural Networks and Learning Systems*, 35:5–22, 2020. URL https://api.semanticscholar.org/CorpusID:220633116. (Cited on 19)

Yiming Li, Yong Jiang, Zhifeng Li, and Shu-Tao Xia. Backdoor learning: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2022. (Cited on 19)

Yuhui Li, Fangyun Wei, Jinjing Zhao, Chao Zhang, and Hongyang Zhang. RAIN: Your Language Models Can Align Themselves without Finetuning. *ArXiv*, abs/2309.07124, 2023d. URL https://api.semanticscholar.org/CorpusID:261705563. (Cited on 24, 28)

Zongjie Li, Chaozheng Wang, Shuai Wang, and Cuiyun Gao. Protecting Intellectual Property of Large Language Model-Based Code Generation APIs via Watermarks. *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, 2023e. URL https://api.semanticscholar.org/CorpusID:265352433. (Cited on 14)

Jiawei Liang, Siyuan Liang, Man Luo, Aishan Liu, Dongchen Han, Ee-Chien Chang, and Xiaochun Cao. VL-Trojan: Multimodal Instruction Backdoor Attacks against Autoregressive Visual Language Models. *ArXiv*, abs/2402.13851, 2024. URL https://api.semanticscholar.org/CorpusID:267770333. (Cited on 6)

Zeyi Liao and Huan Sun. AmpleGCG: Learning a Universal and Transferable Generative Model of Adversarial Suffixes for Jailbreaking Both Open and Closed LLMs. 2024. URL https://api.semanticscholar.org/CorpusID:269043107. (Cited on 4)

Lizhi Lin, Honglin Mu, Zenan Zhai, Minghan Wang, Yuxia Wang, Renxi Wang, Junjie Gao, Yixuan Zhang, Wanxiang Che, Timothy Baldwin, Xudong Han, and Haonan Li. Against The Achilles' Heel: A Survey on Red Teaming for Generative Models. 2024. URL https://api.semanticscholar.org/CorpusID:268820098. (Cited on 6)

Aiwei Liu, Leyi Pan, Yijian Lu, Jingjing Li, Xuming Hu, Lijie Wen, Irwin King, and Philip S. Yu. A survey of text watermarking in the era of large language models. *ACM Comput. Surv.*, 57:47:1–47:36, 2023a. URL https://api.semanticscholar.org/CorpusID:266191530. (Cited on 6, 14)

Guozhi Liu, Weiwei Lin, Tiansheng Huang, Ruichao Mo, Qi Mu, and Li Shen. Targeted vaccine: Safety alignment for large language models against harmful fine-tuning via layer-wise perturbation. *ArXiv*, abs/2410.09760, 2024a. URL https://api.semanticscholar.org/CorpusID:273346246. (Cited on 25, 28)

Jiaxu Liu, Xiangyu Yin, Sihao Wu, Jianhong Wang, Meng Fang, Xinping Yi, and Xiaowei Huang. Tiny Refinements Elicit Resilience: Toward Efficient Prefix-Model Against LLM Red-Teaming. 2024b. URL https://api.semanticscholar.org/CorpusID:269929930. (Cited on 27)

Qin Liu, Fei Wang, Chaowei Xiao, and Muhao Chen. From Shortcuts to Triggers: Backdoor Defense with Denoised PoE. *ArXiv*, abs/2305.14910, 2023b. URL https://api.semanticscholar.org/CorpusID:258866191. (Cited on 28)

Tianqi Liu, Yao Zhao, Rishabh Joshi, Misha Khalman, Mohammad Saleh, Peter J. Liu, and Jialu Liu. Statistical rejection sampling improves preference optimization. *ArXiv*, abs/2309.06657, 2023c. URL https://api.semanticscholar.org/CorpusID:261705578. (Cited on 24)

Tong Liu, Yingjie Zhang, Zhe Zhao, Yinpeng Dong, Guozhu Meng, and Kai Chen. Making Them Ask and Answer: Jailbreaking Large Language Models in Few Queries via Disguise and Reconstruction. *ArXiv*, abs/2402.18104, 2024c. URL https://api.semanticscholar.org/CorpusID:268041845. (Cited on 4)

Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. AutoDAN: Generating Stealthy Jailbreak Prompts on Aligned Large Language Models. *CoRR*, abs/2310.04451, 2023d. doi: 10.48550/ARXIV.2310.04451. URL https://doi.org/10.48550/arXiv.2310.04451. (Cited on 4, 10, 14)

Xiaogeng Liu, Peiran Li, Edward Suh, Yevgeniy Vorobeychik, Zhuoqing Mao, Somesh Jha, Patrick McDaniel, Huan Sun, Bo Li, and Chaowei Xiao. Autodan-turbo: A lifelong agent for strategy self-exploration to jailbreak llms. *CoRR*, abs/2410.05295, 2024d. doi: 10.48550/ARXIV.2410.05295. URL https://doi.org/10.48550/arXiv.2410.05295. (Cited on 14)

Xiaogeng Liu, Zhiyuan Yu, Yizhe Zhang, Ning Zhang, and Chaowei Xiao. Automatic and Universal Prompt Injection Attacks against Large Language Models. 2024e. URL https://api.semanticscholar.org/CorpusID:268296913. (Cited on 4)

Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng, and Yang Liu. Prompt Injection attack against LLM-integrated Applications. *CoRR*, abs/2306.05499, 2023e. doi: 10.48550/ARXIV.2306.05499. URL https://doi.org/10.48550/arXiv.2306.05499. (Cited on 4, 13)

Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, and Yang Liu. Jailbreaking ChatGPT via Prompt Engineering: An Empirical Study. *ArXiv*, abs/2305.13860, 2023f. URL https://api.semanticscholar.org/CorpusID:258841501. (Cited on 4, 11)

Yingqi Liu, Guangyu Shen, Guanhong Tao, Shengwei An, Shiqing Ma, and X. Zhang. Piccolo: Exposing Complex Backdoors in NLP Transformer Models. *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 2025–2042, 2022. URL https://api.semanticscholar.org/CorpusID:248067917. (Cited on 19)

Zichuan Liu, Zefan Wang, Linjie Xu, Jinyu Wang, Lei Song, Tianchun Wang, Chunlin Chen, Wei Cheng, and Jiang Bian. Protecting Your LLMs with Information Bottleneck. *ArXiv*, abs/2404.13968, 2024f. URL https://api.semanticscholar.org/CorpusID:269293591. (Cited on 27)

LLM Agents. LLM Agents. https://www.promptingguide.ai/research/llm-agents, 2023. (Cited on 7, 22)

LLM Security. LLM Security. https://llmsecurity.net/, 2023. (Cited on 22)

Renze Lou, Kai Zhang, and Wenpeng Yin. A Comprehensive Survey on Instruction Following. 2023. URL https://api.semanticscholar.org/CorpusID:257632435. (Cited on 2)

Dong Lu, Tianyu Pang, Chao Du, Qian Liu, Xianjun Yang, and Min Lin. Test-Time Backdoor Attacks on Multimodal Large Language Models. 2024a. URL https://api.semanticscholar.org/CorpusID:267637232. (Cited on 4, 18)

Lin Lu, Hai Yan, Zenghui Yuan, Jiawen Shi, Wenqi Wei, Pin-Yu Chen, and Pan Zhou. AutoJailbreak: Exploring Jailbreak Attacks and Defenses through a Dependency Lens. 2024b. URL https://api.semanticscholar.org/CorpusID:270285580. (Cited on 26)

Huijie Lv, Xiao Wang, Yuan Zhang, Caishuang Huang, Shihan Dou, Junjie Ye, Tao Gui, Qi Zhang, and Xuanjing Huang. CodeChameleon: Personalized Encryption Framework for Jailbreaking Large Language Models. *ArXiv*, abs/2402.16717, 2024. URL https://api.semanticscholar.org/CorpusID:268032340. (Cited on 4, 13)

Chengdong Ma, Ziran Yang, Minquan Gao, Hai Ci, Jun Gao, Xuehai Pan, and Yaodong Yang. Red Teaming Game: A Game-Theoretic Framework for Red Teaming Language Models. *ArXiv*, abs/2310.00322, 2023. URL https://api.semanticscholar.org/CorpusID:263334034. (Cited on 28)

Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables. *ArXiv*, abs/1611.00712, 2016. URL https://api.semanticscholar.org/CorpusID:14307651. (Cited on 21)

Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards Deep Learning Models Resistant to Adversarial Attacks. *ArXiv*, abs/1706.06083, 2017. URL https://api.semanticscholar.org/CorpusID:3488815. (Cited on 24)

Vahid Majdinasab, Michael Joshua Bishop, Shawn Rasheed, Arghavan Moradi Dakhel, Amjed Tahir, and Foutse Khomh. Assessing the Security of GitHub Copilot Generated Code - A Targeted Replication Study. *ArXiv*, abs/2311.11177, 2023. URL https://api.semanticscholar.org/CorpusID:265295550. (Cited on 1)

Neal Mangaokar, Ashish Hooda, Jihye Choi, Shreyas Chandrashekaran, Kassem Fawaz, Somesh Jha, and Atul Prakash. PRP: Propagating Universal Perturbations to Attack Large Language Model Guard-Rails. *ArXiv*, abs/2402.15911, 2024. URL https://api.semanticscholar.org/CorpusID:267938152. (Cited on 22)

Todor Markov, Chong Zhang, Sandhini Agarwal, Florentine Eloundou Nekoul, Theodore Lee, Steven Adler, Angela Jiang, and Lilian Weng. A Holistic Approach to Undesired Content Detection in the Real World. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(12):15009–15018, Jun. 2023. doi: 10.1609/aaai.v37i12.26752. URL https://ojs.aaai.org/index.php/AAAI/article/view/26752. (Cited on 22)

Thomas Marshall, Adam Scherlis, and Nora Belrose. Refusal in llms is an affine function. *CoRR*, abs/2411.09003, 2024. doi: 10.48550/ARXIV.2411.09003. URL https://doi.org/10.48550/arXiv.2411.09003. (Cited on 18)

Justus Mattern, Fatemehsadat Mireshghallah, Zhijing Jin, Bernhard Schoelkopf, Mrinmaya Sachan, and Taylor Berg-Kirkpatrick. Membership inference attacks against language models via neighbourhood comparison. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 11330–11343, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.719. URL https://aclanthology.org/2023.findings-acl.719. (Cited on 4)

Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, David Forsyth, and Dan Hendrycks. HarmBench: A Standardized Evaluation Framework for Automated Red Teaming and Robust Refusal. *ArXiv*, abs/2402.04249, 2024. URL https://api.semanticscholar.org/CorpusID:267499790. (Cited on 3, 29)

Matthieu Meeus, Shubham Jain, Marek Rei, and Yves-Alexandre de Montjoye. Did the neurons read your book? document-level membership inference for large language models. *arXiv preprint arXiv:2310.15007*, 2023. (Cited on 15)

Ninareh Mehrabi, Palash Goyal, Christophe Dupuy, Qian Hu, Shalini Ghosh, Richard S. Zemel, Kai-Wei Chang, Aram Galstyan, and Rahul Gupta. FLIRT: Feedback Loop In-context Red Teaming. *CoRR*, abs/2308.04265, 2023a. doi: 10.48550/ARXIV.2308.04265. URL https://doi.org/10.48550/arXiv.2308.04265. (Cited on 4, 13)

Ninareh Mehrabi, Palash Goyal, Anil Ramakrishna, J. Dhamala, Shalini Ghosh, Richard Zemel, Kai-Wei Chang, A. G. Galstyan, and Rahul Gupta. JAB: Joint Adversarial Prompting and Belief Augmentation. *ArXiv*, abs/2311.09473, 2023b. URL https://api.semanticscholar.org/CorpusID:265220687. (Cited on 28)

Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. Tree of Attacks: Jailbreaking Black-Box LLMs Automatically. *ArXiv*, abs/2312.02119, 2023. URL https://api.semanticscholar.org/CorpusID:265609901. (Cited on 4)

Kai Mei, Zelong Li, Shuyuan Xu, Ruosong Ye, Yingqiang Ge, and Yongfeng Zhang. LLM Agent Operating System. 2024a. URL https://api.semanticscholar.org/CorpusID:268681490. (Cited on 22)

Lingrui Mei, Shenghua Liu, Yiwei Wang, Baolong Bi, Jiayi Mao, and Xueqi Cheng. "Not Aligned"is Not"Malicious": Being Careful about Hallucinations of Large Language Models' Jailbreak. 2024b. URL https://api.semanticscholar.org/CorpusID:270559880. (Cited on 29, 31)

Peter Mell, Karen Scarfone, and Sasha Romanosky. A complete guide to the common vulnerability scoring system version 2.0, 2007-07-30 2007. URL https://tsapps.nist.gov/publication/get_pdf.cfm?pub_id=51198. (Cited on 31)

Meta. Overview of Meta AI safety policies prepared for the UK AI Safety Summit. https://transparency.fb.com/en-gb/policies/ai-safety-policies-for-safety-summit/, 2023. (Cited on 2)

Microsoft. The Python Risk Identification Tool for generative AI (PyRIT). https://github.com/Azure/PyRIT, 2023. (Cited on 29)

Fatemehsadat Mireshghallah, Kartik Goyal, Archit Uniyal, Taylor Berg-Kirkpatrick, and Reza Shokri. Quantifying privacy risks of masked language models using membership inference attacks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 8332–8347, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL https://aclanthology.org/2022.emnlp-main.570. (Cited on 15)

Lingbo Mo, Zeyi Liao, Boyuan Zheng, Yu Su, Chaowei Xiao, and Huan Sun. A Trembling House of Cards? Mapping Adversarial Attacks against Language Agents. *ArXiv*, abs/2402.10196, 2024a. URL https://api.semanticscholar.org/CorpusID:267682286. (Cited on 22)

Yichuan Mo, Yuji Wang, Zeming Wei, and Yisen Wang. Studious Bob Fight Back Against Jailbreaking via Prompt Adversarial Tuning. *ArXiv*, abs/2402.06255, 2024b. URL https://api.semanticscholar.org/CorpusID:267617138. (Cited on 22, 27)

Model Denial of Service. Model Denial of Service. https://wiki.hego.tech/owasp/owasp-llm-top-10-v1.0/llm04-model-denial-of-service, 2023. (Cited on 9)

John Morris, Volodymyr Kuleshov, Vitaly Shmatikov, and Alexander Rush. Text Embeddings Reveal (Almost) As Much As Text. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12448–12460, Singapore, December 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.765. URL https://aclanthology.org/2023.emnlp-main.765. (Cited on 16, 26)

John X. Morris, Volodymyr Kuleshov, Vitaly Shmatikov, and Alexander M. Rush. Text embeddings reveal (almost) as much as text. In *Conference on Empirical Methods in Natural Language Processing*, 2023b. URL https://api.semanticscholar.org/CorpusID:263829206. (Cited on 4)

John X. Morris, Wenting Zhao, Justin T. Chiu, Vitaly Shmatikov, and Alexander M. Rush. Language Model Inversion. *CoRR*, abs/2311.13647, 2023c. doi: 10.48550/ARXIV.2311.13647. URL https://doi.org/10.48550/arXiv.2311.13647. (Cited on 4, 15, 16)

Maximilian Mozes, Xuanli He, Bennett Kleinberg, and Lewis D. Griffin. Use of LLMs for Illicit Purposes: Threats, Prevention Measures, and Vulnerabilities. *ArXiv*, abs/2308.12833, 2023. URL https://api.semanticscholar.org/CorpusID:261101245. (Cited on 22)

Arif Ali Mughal. The art of cybersecurity: Defense in depth strategy for robust protection. *International Journal of Intelligent Automation and Computing*, 1(1):1–20, Mar. 2018. URL https://research.tensorgate.org/index.php/IJIAC/article/view/19. (Cited on 30, 31)

Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland, Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Andrea Michi, Marco Selvi, Sertan Girgin, Nikola Momchev, Olivier Bachem, Daniel Jaymin Mankowitz, Doina Precup, and Bilal Piot. Nash learning from human feedback. *ArXiv*, abs/2312.00886, 2023. URL https://api.semanticscholar.org/CorpusID:265609682. (Cited on 24)

Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable Extraction of Training Data from (Production) Language Models. *ArXiv*, abs/2311.17035, 2023. URL https://api.semanticscholar.org/CorpusID:265466445. (Cited on 4, 15)

NIST CVE. NIST CVE. https://nvd.nist.gov/vuln, 2022. (Cited on 3)

Zhenxing Niu, Haodong Ren, Xinbo Gao, Gang Hua, and Rong Jin. Jailbreaking Attack against Multimodal Large Language Model. *ArXiv*, abs/2402.02309, 2024. URL https://api.semanticscholar.org/CorpusID:267413270. (Cited on 6)

Munachiso Nwadike, Takumi Miyawaki, Esha Sarkar, Michail Maniatakos, and Farah Shamout. Explainability matters: Backdoor attacks on medical imaging. *arXiv preprint arXiv:2101.00008*, 2020. (Cited on 6)

Myung Gyo Oh, Leo Hyun Park, Jaeuk Kim, Jaewoo Park, and Taekyoung Kwon. Membership Inference Attacks With Token-Level Deduplication on Korean Language Models. *IEEE Access*, 11:10207–10217, 2023. (Cited on 15)

Seong Joon Oh, Maximilian Augustin, Mario Fritz, and Bernt Schiele. Towards Reverse-Engineering Black-Box Neural Networks. In *International Conference on Learning Representations*, 2017. URL https://api.semanticscholar.org/CorpusID:3278569. (Cited on 14)

OpenAI. Introducing ChatGPT. https://openai.com/blog/chatgpt, 2022. (Cited on 7, 14)

OpenAI. GPT-4 Technical Report. 2023. URL https://api.semanticscholar.org/CorpusID:257532815. (Cited on 1)

OpenAI. Function Calling OpenAI. https://platform.openai.com/docs/guides/function-calling, 2023. (Cited on 12)

OpenAI. Using logit bias to define token probability. https://help.openai.com/en/articles/5247780-using-logit-bias-to-define-token-probability, 2023. (Cited on 16)

OpenAI. OpenAI Red Teaming Network. https://openai.com/blog/red-teaming-network, 2023. (Cited on 30)

OpenAI. GPT4 Fine Tuning. https://platform.openai.com/docs/guides/fine-tuning, 2024a. (Cited on 20)

OpenAI. OpenAI Moderation Endpoint. platform.openai.com/docs/guides/moderation/overview, 2024b. (Cited on 3, 22, 27)

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training Language Models to Follow Instructions with Human Feedback. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022a. URL https://openreview.net/forum?id=TG8KACxEON. (Cited on 2)

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022b. (Cited on 2, 24)

OWASP. 2025 Top 10 Risk Mitigations for LLMs and Gen AI Apps. https://genaisecurityproject.com/llm-top-10/, 2025. (Cited on 3)

Mustafa Safa Ozdayi, Charith S. Peris, Jack G. M. FitzGerald, Christophe Dupuy, Jimit Majmudar, Haidar Khan, Rahil Parikh, and Rahul Gupta. Controlling the Extraction of Memorized Data from Large Language Models via Prompt-Tuning. In *Annual Meeting of the Association for Computational Linguistics*, 2023. URL https://api.semanticscholar.org/CorpusID:258823013. (Cited on 4, 25)

Nicolas Papernot, Patrick Mcdaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical Black-Box Attacks against Machine Learning. *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*, 2016. URL https://api.semanticscholar.org/CorpusID:1090603. (Cited on 14)

Nicolas Papernot, Patrick Mcdaniel, Arunesh Sinha, and Michael P. Wellman. SoK: Security and Privacy in Machine Learning. *2018 IEEE European Symposium on Security and Privacy (EuroS&P)*, pp. 399–414, 2018. URL https://api.semanticscholar.org/CorpusID:44237208. (Cited on 9, 10, 11)

Dario Pasquini, Martin Strohmeier, and Carmela Troncoso. Neural Exec: Learning (and Learning from) Execution Triggers for Prompt Injection Attacks. *ArXiv*, abs/2403.03792, 2024. URL https://api.semanticscholar.org/CorpusID:268253024. (Cited on 4)

Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large Language Model Connected with Massive APIs. *ArXiv*, abs/2305.15334, 2023. URL https://api.semanticscholar.org/CorpusID:258865184. (Cited on 7)

Sonali Pattnaik, Rohan Karan, Srijan Kumar, and Clémentine Fourrier. Introducing the Chatbot Guardrails Arena. https://huggingface.co/blog/arena-lighthouz, 2024. (Cited on 29)

Anselm Paulus, Arman Zharmagambetov, Chuan Guo, Brandon Amos, and Yuandong Tian. AdvPrompter: Fast Adaptive Adversarial Prompting for LLMs. *ArXiv*, abs/2404.16873, 2024. URL https://api.semanticscholar.org/CorpusID:269430799. (Cited on 14)

Will Pearce and Joseph Lucas. NVIDIA AI Red Team: An Introduction. https://developer.nvidia.com/blog/nvidia-ai-red-team-an-introduction/, 2023. (Cited on 2)

Kellin Pelrine, Mohammad Taufeeque, Michal Zajac, Euan McLean, and Adam Gleave. Exploiting Novel GPT-4 APIs. 2023. URL https://api.semanticscholar.org/CorpusID:266521205. (Cited on 4, 12)

Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra-Aimée Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The RefinedWeb Dataset for Falcon LLM: Outperforming Curated Corpora with Web Data, and Web Data Only. *ArXiv*, abs/2306.01116, 2023. URL https://api.semanticscholar.org/CorpusID:259063761. (Cited on 16)

Sheng-Hsuan Peng, Pin-Yu Chen, Matthew Hull, and Duen Horng Chau. Navigating the Safety Landscape: Measuring Risks in Finetuning Large Language Models. *ArXiv*, abs/2405.17374, 2024a. URL https://api.semanticscholar.org/CorpusID:270063157. (Cited on 20)

Shengyun Peng, Pin-Yu Chen, Matthew Hull, and Duen Horng Chau. Navigating the safety landscape: Measuring risks in finetuning large language models. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024b. URL http://papers.nips.cc/paper_files/paper/2024/hash/ada93fa6643735f294be51dc31eebbd4-Abstract-Conference.html. (Cited on 24)

Yu Peng, Zewen Long, Fangming Dong, Congyi Li, Shu Wu, and Kai Chen. Playing language game with llms leads to jailbreaking. *CoRR*, abs/2411.12762, 2024c. doi: 10.48550/ARXIV.2411.12762. URL https://doi.org/10.48550/arXiv.2411.12762. (Cited on 12)

Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 3419–3448, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL https://aclanthology.org/2022.emnlp-main.225. (Cited on 1, 4, 12, 21, 28, 29)

Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson Elhage, Nicholas Joseph, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. Discovering language model behaviors with model-written evaluations. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 13387–13434, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.847. URL https://aclanthology.org/2023.findings-acl.847. (Cited on 1)

Perspective API. Perspective API Toxicity Categories. https://developers.perspectiveapi.com/s/about-the-api-attributes-and-languages, 2024. (Cited on 3)

Willem Pienaar and Shahram Anver. Rebuff: Detecting Prompt Injection Attacks. https://blog.langchain.dev/rebuff/, 2023. (Cited on 26)

Jeff Price. Red Alert: Are you ready for every airport manager's worst nightmare? *Airport Magazine*, 16, 2004. URL https://api.semanticscholar.org/CorpusID:106638276. (Cited on 1)

Fanchao Qi, Yangyi Chen, Mukai Li, Zhiyuan Liu, and Maosong Sun. ONION: A Simple and Effective Defense Against Textual Backdoor Attacks. *ArXiv*, abs/2011.10369, 2020. URL https://api.semanticscholar.org/CorpusID:227118606. (Cited on 23)

Fanchao Qi, Yangyi Chen, Xurui Zhang, Mukai Li, Zhiyuan Liu, and Maosong Sun. Mind the Style of Text! Adversarial and Backdoor Attacks Based on Text Style Transfer. *ArXiv*, abs/2110.07139, 2021a. URL https://api.semanticscholar.org/CorpusID:238857078. (Cited on 19)

Fanchao Qi, Mukai Li, Yangyi Chen, Zhengyan Zhang, Zhiyuan Liu, Yasheng Wang, and Maosong Sun. Hidden Killer: Invisible Textual Backdoor Attacks with Syntactic Trigger. In *Annual Meeting of the Association for Computational Linguistics*, 2021b. URL https://api.semanticscholar.org/CorpusID:235196099. (Cited on 19)

Fanchao Qi, Yuan Yao, Sophia Xu, Zhiyuan Liu, and Maosong Sun. Turn the Combination Lock: Learnable Textual Backdoor Attacks via Word Substitution. In *Annual Meeting of the Association for Computational Linguistics*, 2021c. URL https://api.semanticscholar.org/CorpusID:235417102. (Cited on 19)

Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning Aligned Language Models Compromises Safety, Even When Users Do Not Intend To! *CoRR*, abs/2310.03693, 2023. doi: 10.48550/ARXIV.2310.03693. URL https://doi.org/10.48550/arXiv.2310.03693. (Cited on 4, 20)

Xiangyu Qi, Yangsibo Huang, Yi Zeng, Edoardo Debenedetti, Jonas Geiping, Luxi He, Kaixuan Huang, Udari Madhushani, Vikash Sehwag, Weijia Shi, Boyi Wei, Tinghao Xie, Danqi Chen, Pin-Yu Chen, Jeffrey Ding, Ruoxi Jia, Jiaqi Ma, Arvind Narayanan, Weijie J. Su, Mengdi Wang, Chaowei Xiao, Bo Li, Dawn Song, Peter Henderson, and Prateek Mittal. AI risk management should incorporate both safety and security. *CoRR*, abs/2405.19524, 2024a. doi: 10.48550/ARXIV.2405.19524. URL https://doi.org/10.48550/arXiv.2405.19524. (Cited on 5)

Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami, Prateek Mittal, and Peter Henderson. Safety Alignment Should Be Made More Than Just a Few Tokens Deep. 2024b. URL https://api.semanticscholar.org/CorpusID:270371778. (Cited on 27)

Yao Qiang, Xiangyu Zhou, and Dongxiao Zhu. Hijacking Large Language Models via Adversarial In-Context Learning. *ArXiv*, abs/2311.09948, 2023. URL https://api.semanticscholar.org/CorpusID:265221164. (Cited on 4, 21)

Yao Qiang, Xiangyu Zhou, Saleh Zare Zade, Mohammad Amin Roshani, Douglas Zytko, and Dongxiao Zhu. Learning to Poison Large Language Models During Instruction Tuning. 2024. URL https://api.semanticscholar.org/CorpusID:267770200. (Cited on 4)

Baha Rababah, Shang Tommy Wu, Matthew Kwiatkowski, Carson K. Leung, and Cuneyt Gurcan Akcora. Sok: Prompt hacking of large language models. In *2024 IEEE International Conference on Big Data (BigData)*, pp. 5392–5401, 2024. doi: 10.1109/BigData62323.2024.10825103. (Cited on 12)

Bhaktipriya Radharapu, Kevin Robinson, Lora Aroyo, and Preethi Lahoti. AART: AI-Assisted Red-Teaming with Diverse Data Generation for New LLM-powered Applications. In Mingxuan Wang and Imed Zitouni (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: EMNLP 2023 - Industry Track, Singapore, December 6-10, 2023*, pp. 380–395. Association for Computational Linguistics, 2023. URL https://aclanthology.org/2023.emnlp-industry.37. (Cited on 4, 13)

Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023. (Cited on 2)

Parijat Rai, Saumil Sood, Vijay Krishna Madisetti, and Arshdeep Bahga. GUARDIAN: A Multi-Tiered Defense Architecture for Thwarting Prompt Injection Attacks on LLMs. *Journal of Software Engineering and Applications*, 2024. URL https://api.semanticscholar.org/CorpusID:267125160. (Cited on 26)

Vyas Raina, Adian Liusie, and Mark Gales. Is LLM-as-a-Judge Robust? Investigating Universal Adversarial Attacks on Zero-shot LLM Assessment. 2024. URL https://api.semanticscholar.org/CorpusID:267770121. (Cited on 29)

Swaroop Indra Ramaswamy, Om Thakkar, Rajiv Mathews, Galen Andrew, H. B. McMahan, and Franccoise Beaufays. Training Production Language Models without Memorizing User Data. *ArXiv*, abs/2009.10031, 2020. URL https://api.semanticscholar.org/CorpusID:221819648. (Cited on 25)

Javier Rando and Florian Tramèr. Universal Jailbreak Backdoors from Poisoned Human Feedback. *ArXiv*, abs/2311.14455, 2023. URL https://api.semanticscholar.org/CorpusID:265445004. (Cited on 4, 19)

Javier Rando and Florian Tramèr. RLHF Trojan Competition. https://github.com/ethz-spylab/rlhf_trojan_competition, 2024a. (Cited on 2, 19, 20)

Javier Rando and Florian Tramèr. RLHF Trojan Competition. twitter.com/javirandor/status/1762828917872730324, 2024b. (Cited on 2, 20)

Abhinav Rao, Sachin Vashistha, Atharva Naik, Somak Aditya, and Monojit Choudhury. Tricking LLMs into Disobedience: Understanding, Analyzing, and Preventing Jailbreaks. *arXiv preprint arXiv:2305.14965*, 2023. (Cited on 4)

James Reason. *Human error*. Cambridge University Press, 1990. (Cited on 25)

Traian Rebedea, Razvan Laurentiu Dinu, Makesh Narsimhan Sreedhar, Christopher Parisien, and Jonathan Cohen. NeMo Guardrails: A Toolkit for Controllable and Safe LLM Applications with Programmable Rails. In *Conference on Empirical Methods in Natural Language Processing*, 2023. URL https://api.semanticscholar.org/CorpusID:264146531. (Cited on 22, 27)

Red Team. Red Team Wikipedia. https://en.wikipedia.org/wiki/Red_team. (Cited on 1)

Richard Ren, Steven Basart, Adam Khoja, Alice Gatti, Long Phan, Xuwang Yin, Mantas Mazeika, Alexander Pan, Gabriel Mukobi, Ryan H. Kim, Stephen Fitz, and Dan Hendrycks. Safetywashing: Do ai safety benchmarks actually measure safety progress?, 2024. URL https://arxiv.org/abs/2407.21792. (Cited on 29)

Maria Rigaki and Sebastian Garcia. A Survey of Privacy Attacks in Machine Learning. *ACM Comput. Surv.*, 56(4), nov 2023. ISSN 0360-0300. doi: 10.1145/3624010. URL https://doi.org/10.1145/3624010. (Cited on 15)

Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. SmoothLLM: Defending Large Language Models Against Jailbreaking Attacks. *ArXiv*, abs/2310.03684, 2023. URL https://api.semanticscholar.org/CorpusID:263671542. (Cited on 22, 27)

Christophe Ropers, David Dale, Prangthip Hansanti, Gabriel Mejia Gonzalez, Ivan Evtimov, Corinne Wong, Christophe Touret, Kristina Pereyra, Seohyun Sonia Kim, Cristian Canton-Ferrer, Pierre Andrews, and Marta Ruiz Costa-jussà. Towards Red Teaming in Multimodal and Multilingual Translation. *ArXiv*, abs/2401.16247, 2024. URL https://api.semanticscholar.org/CorpusID:267311645. (Cited on 29)

Domenic Rosati, Jan Wehner, Kai Williams, Lukasz Bartoszcze, David Atanasov, Robie Gonzales, Subhabrata Majumdar, Carsten Maple, Hassan Sajjad, and Frank Rudzicz. Representation noising effectively prevents harmful fine-tuning on LLMs. 2024a. URL https://api.semanticscholar.org/CorpusID:269982864. (Cited on 25, 27)

Domenic Rosati, Jan Wehner, Kai Williams, Lukasz Bartoszcze, Jan Batzner, Hassan Sajjad, and Frank Rudzicz. Immunization against harmful fine-tuning attacks. 2024b. URL https://api.semanticscholar.org/CorpusID:268032044. (Cited on 25, 27)

Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. XSTest: A Test Suite for Identifying Exaggerated Safety Behaviours in Large Language Models. *ArXiv*, abs/2308.01263, 2023. URL https://api.semanticscholar.org/CorpusID:260378842. (Cited on 24, 30, 31)

Jessica Rumbelow and Matthew Watkins. SolidGoldMagikarp (plus, prompt generation). lesswrong.com/solidgoldmagikarp-plus-prompt-generation, 2023. (Cited on 12)

Mark Russinovich. Mitigating Skeleton Key, a new type of generative AI jailbreak technique. https://www.microsoft.com/en-us/security/blog/2024/06/26/mitigating-skeleton-key-a-new-type-of-generative-ai-jailbreak-technique/, 2024. (Cited on 4, 11)

Shiva Omrani Sabbaghi, Robert Wolfe, and Aylin Caliskan. Evaluating Biased Attitude Associations of Language Models in an Intersectional Context. *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, 2023. URL https://api.semanticscholar.org/CorpusID:259375792. (Cited on 1)

Vinu Sankar Sadasivan, Shoumik Saha, Gaurang Sriramanan, Priyatham Kattakinda, Atoosa Malemir Chegini, and Soheil Feizi. Fast Adversarial Attacks on Language Models In One GPU Minute. *ArXiv*, abs/2402.15570, 2024. URL https://api.semanticscholar.org/CorpusID:267938703. (Cited on 4, 18)

Mikayel Samvelyan, Sharath Chandra Raparthy, Andrei Lupu, Eric Hambro, Aram H. Markosyan, Manish Bhatt, Yuning Mao, Minqi Jiang, Jack Parker-Holder, Jakob Foerster, Tim Rocktaschel, and Roberta Raileanu. Rainbow Teaming: Open-Ended Generation of Diverse Adversarial Prompts. *ArXiv*, abs/2402.16822, 2024. URL https://api.semanticscholar.org/CorpusID:268031888. (Cited on 28)

Omar Sanseviero. Pickle Scanning, Safetensors, Social validation features. twitter.com/osanseviero/status/1763331704146583806, 2024. (Cited on 6)

Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. Whose Opinions Do Language Models Reflect? *ArXiv*, abs/2303.17548, 2023. URL https://api.semanticscholar.org/CorpusID:257834040. (Cited on 1)

William Saunders, Catherine Yeh, Jeff Wu, Steven Bills, Long Ouyang, Jonathan Ward, and Jan Leike. Self-critiquing models for assisting human evaluators. *CoRR*, abs/2206.05802, 2022. doi: 10.48550/ARXIV.2206.05802. URL https://doi.org/10.48550/arXiv.2206.05802. (Cited on 24)

Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language Models Can Teach Themselves to Use Tools. *ArXiv*, abs/2302.04761, 2023. URL https://api.semanticscholar.org/CorpusID:256697342. (Cited on 7)

Jonas Schuett, Noemi Dreksler, Markus Anderljung, David McCaffary, Lennart Heim, Emma Bluemke, and Ben Garfinkel. Towards best practices in AGI safety and governance: A survey of expert opinion. *ArXiv*, abs/2305.07153, 2023. URL https://api.semanticscholar.org/CorpusID:258676345. (Cited on 6)

Sander Schulhoff, Jeremy Pinto, Anaum Khan, Louis-François Bouchard, Chenglei Si, Svetlina Anati, Valen Tagliabue, Anson Kost, Christopher Carnahan, and Jordan Boyd-Graber. Ignore this title and HackAPrompt: Exposing systemic vulnerabilities of LLMs through a global prompt hacking competition. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 4945–4977, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.302. URL https://aclanthology.org/2023.emnlp-main.302. (Cited on 4, 6, 11, 12)

Barry Schwartz. Bing Webmaster Guidelines. https://searchengineland.com/prompt-injection-added-to-bing-webmaster-guidelines-443788, 2024. (Cited on 18)

Leo Schwinn, David Dobre, Sophie Xhonneux, Gauthier Gidel, and Stephan Günnemann. Soft Prompt Threats: Attacking Safety Alignment and Unlearning in Open-Source LLMs through the Embedding Space. *ArXiv*, abs/2402.09063, 2024. URL https://api.semanticscholar.org/CorpusID:267657556. (Cited on 4)

Rico Sennrich, Jannis Vamvas, and Alireza Mohammadshahi. Mitigating Hallucinations and Off-target Machine Translation with Source-Contrastive and Language-Contrastive Decoding. *ArXiv*, abs/2309.07098, 2023. URL https://api.semanticscholar.org/CorpusID:261705560. (Cited on 24)

Erfan Shayegani, Md Abdullah Al Mamun, Yu Fu, Pedram Zaree, Yue Dong, and Nael B. Abu-Ghazaleh. Survey of Vulnerabilities in Large Language Models Revealed by Adversarial Attacks. *ArXiv*, abs/2310.10844, 2023. URL https://api.semanticscholar.org/CorpusID:264172191. (Cited on 6)

Guangyu Shen, Siyuan Cheng, Kai xian Zhang, Guanhong Tao, Shengwei An, Lu Yan, Zhuo Zhang, Shiqing Ma, and Xiangyu Zhang. Rapid Optimization for Jailbreaking LLMs via Subconscious Exploitation and Echopraxia. *ArXiv*, abs/2402.05467, 2024. URL https://api.semanticscholar.org/CorpusID:267547566. (Cited on 4)

Lujia Shen, Shouling Ji, Xuhong Zhang, Jinfeng Li, Jing Chen, Jie Shi, Chengfang Fang, Jianwei Yin, and Ting Wang. Backdoor Pre-trained Models Can Transfer to All. *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, 2021. URL https://api.semanticscholar.org/CorpusID:240354696. (Cited on 19)

Sheng Shen, Le Hou, Yan-Quan Zhou, Nan Du, S. Longpre, Jason Wei, Hyung Won Chung, Barret Zoph, William Fedus, Xinyun Chen, Tu Vu, Yuexin Wu, Wuyang Chen, Albert Webson, Yunxuan Li, Vincent Zhao, Hongkun Yu, Kurt Keutzer, Trevor Darrell, and Denny Zhou. Mixture-of-Experts Meets Instruction Tuning:A Winning Combination for Large Language Models. 2023a. URL https://api.semanticscholar.org/CorpusID:259342096. (Cited on 17)

Tianhao Shen, Renren Jin, Yufei Huang, Chuang Liu, Weilong Dong, Zishan Guo, Xinwei Wu, Yan Liu, and Deyi Xiong. Large Language Model Alignment: A Survey. *ArXiv*, abs/2309.15025, 2023b. URL https://api.semanticscholar.org/CorpusID:262824801. (Cited on 2)

Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. "Do Anything Now": Characterizing and Evaluating In-The-Wild Jailbreak Prompts on Large Language Models. *CoRR abs/2308.03825*, 2023c. (Cited on 1)

Xinyue Shen, Zeyuan Johnson Chen, Michael Backes, Yun Shen, and Yang Zhang. "Do Anything Now": Characterizing and Evaluating In-The-Wild Jailbreak Prompts on Large Language Models. *ArXiv*, abs/2308.03825, 2023d. URL https://api.semanticscholar.org/CorpusID:260704242. (Cited on 4, 10, 11)

Xuan Sheng, Zhaoyang Han, Piji Li, and Xiangmao Chang. A survey on backdoor attack and defense in natural language processing. In *2022 IEEE 22nd International Conference on Software Quality, Reliability and Security (QRS)*, pp. 809–820. IEEE, 2022. (Cited on 19)

Xuan Sheng, Zhicheng Li, Zhaoyang Han, Xiangmao Chang, and Piji Li. Punctuation Matters! Stealthy Backdoor Attack for Language Models. In Fei Liu, Nan Duan, Qingting Xu, and Yu Hong (eds.), *Natural Language Processing and Chinese Computing - 12th National CCF Conference, NLPCC 2023, Foshan, China, October 12-15, 2023, Proceedings, Part I*, volume 14302 of *Lecture Notes in Computer Science*, pp. 524–536. Springer, 2023. doi: 10.1007/978-3-031-44693-1\\_41. URL https://doi.org/10.1007/978-3-031-44693-1_41. (Cited on 19)

Toby Shevlane, Sebastian Farquhar, Ben Garfinkel, Mary Phuong, Jess Whittlestone, Jade Leung, Daniel Kokotajlo, Nahema Marchal, Markus Anderljung, Noam Kolt, Lewis Ho, Divya Siddarth, Shahar Avin, William T. Hawkins, Been Kim, Iason Gabriel, Vijay Bolina, Jack Clark, Yoshua Bengio, Paul Francis Christiano, and Allan Dafoe. Model evaluation for extreme risks. *ArXiv*, abs/2305.15324, 2023. URL https://api.semanticscholar.org/CorpusID:258865507. (Cited on 6)

R. Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership Inference Attacks Against Machine Learning Models. *2017 IEEE Symposium on Security and Privacy (SP)*, pp. 3–18, 2016. URL https://api.semanticscholar.org/CorpusID:10488675. (Cited on 14)

Manli Shu, Jiong Wang, Chen Zhu, Jonas Geiping, Chaowei Xiao, and Tom Goldstein. On the Exploitability of Instruction Tuning. *ArXiv*, abs/2306.17194, 2023. URL https://api.semanticscholar.org/CorpusID:259309096. (Cited on 4)

Aradhana Sinha, Ananth Balashankar, Ahmad Beirami, Thi Avrahami, Jilin Chen, and Alex Beutel. Break it, Imitate it, Fix it: Robustness by Generating Human-Like Attacks. *ArXiv*, abs/2310.16955, 2023. URL https://api.semanticscholar.org/CorpusID:264490789. (Cited on 28)

Chawin Sitawarin, Norman Mu, David Wagner, and Alexandre Araujo. PAL: Proxy-Guided Black-Box Attack on Large Language Models. *ArXiv*, abs/2402.09674, 2024. URL https://api.semanticscholar.org/CorpusID:267682038. (Cited on 4)

Ram Shankar Siva Kumar. A Few Useful Lessons about AI Red Teaming. https://youtu.be/UKj5jj6fD_c?t=2025, 2023a. (Cited on 6)

Ram Shankar Siva Kumar. Microsoft AI Red Team building future of safer AI. https://www.microsoft.com/en-us/security/blog/2023/08/07/microsoft-ai-red-team-building-future-of-safer-ai/, 2023b. (Cited on 2)

Victoria Smith, Ali Shahin Shamsabadi, Carolyn Ashurst, and Adrian Weller. Identifying and Mitigating Privacy Risks Stemming from Language Models: A Survey. *ArXiv*, abs/2310.01424, 2023. URL https://api.semanticscholar.org/CorpusID:263608702. (Cited on 25)

Saleh Soltan, Shankar Ananthakrishnan, Jack FitzGerald, Rahul Gupta, Wael Hamza, Haidar Khan, Charith Peris, Stephen Rawls, Andy Rosenbaum, Anna Rumshisky, Chandana Satya Prakash, Mukund Sridhar, Fabian Triefenbach, Apurv Verma, Gökhan Tür, and Prem Natarajan. AlexaTM 20B: Few-Shot Learning Using a Large-Scale Multilingual Seq2Seq Model. *CoRR*, abs/2208.01448, 2022. doi: 10.48550/ARXIV.2208.01448. URL https://doi.org/10.48550/arXiv.2208.01448. (Cited on 7)

Gowthami Somepalli, Vasu Singla, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Diffusion Art or Digital Forgery? Investigating Data Replication in Diffusion Models. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6048–6058, 2022. URL https://api.semanticscholar.org/CorpusID:254366634. (Cited on 14)

Congzheng Song and Ananth Raghunathan. Information Leakage in Embedding Models. *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 2020. URL https://api.semanticscholar.org/CorpusID:214743021. (Cited on 16)

Shuang Song, Kamalika Chaudhuri, and Anand D. Sarwate. Stochastic gradient descent with differentially private updates. *2013 IEEE Global Conference on Signal and Information Processing*, pp. 245–248, 2013. URL https://api.semanticscholar.org/CorpusID:8085841. (Cited on 25)

Murugiah Souppaya, Karen Scarfone, and Donna Dodson. Secure software development framework (ssdf) version 1.1. *NIST Special Publication*, 800:218, 2022. (Cited on 31)

Logan Stapleton, Jordan Taylor, Sarah Fox, Tongshuang Sherry Wu, and Haiyi Zhu. Seeing Seeds Beyond Weeds: Green Teaming Generative AI for Beneficial Uses. *ArXiv*, abs/2306.03097, 2023. URL https://api.semanticscholar.org/CorpusID:259088586. (Cited on 3, 29)

Asa Cooper Stickland, Alexander Lyzhov, Jacob Pfau, Salsabila Mahdi, and Samuel R. Bowman. Steering Without Side Effects: Improving Post-Deployment Control of Language Models. 2024. URL https://api.semanticscholar.org/CorpusID:270703306. (Cited on 27)

Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zheng Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chun-Yan Li, Eric P. Xing, Furong Huang, Haodong Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Sekhar Jana, Tian-Xiang Chen, Tianming Liu, Tianying Zhou, William Wang, Xiang Li, Xiang-Yu Zhang, Xiao Wang, Xingyao Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, and Yue Zhao. TrustLLM: Trustworthiness in Large Language Models. *ArXiv*, abs/2401.05561, 2024. URL https://api.semanticscholar.org/CorpusID:266933236. (Cited on 5)

Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. A Simple and Effective Pruning Approach for Large Language Models. *ArXiv*, abs/2306.11695, 2023. URL https://api.semanticscholar.org/CorpusID:259203115. (Cited on 27)

Xuchen Suo. Signed-Prompt: A New Approach to Prevent Prompt Injection Attacks Against LLM-Integrated Applications. *ArXiv*, abs/2401.07612, 2024. URL https://api.semanticscholar.org/CorpusID:266999840. (Cited on 27)

Alex Tamkin, Miles Brundage, Jack Clark, and Deep Ganguli. Understanding the capabilities, limitations, and societal impact of large language models. *arXiv preprint arXiv:2102.02503*, 2021. (Cited on 5)

Liyan Tang, Z. Sun, Betina R.S. Idnay, Jordan G. Nestor, Amani Soroush, Pablo Adolfo Elias, Z. P. Xu, Y. Ding, Greg Durrett, Justin F. Rousseau, Cindy Weng, and Yalei Peng. Evaluating large language models on medical evidence summarization. *NPJ Digital Medicine*, 6, 2023. URL https://api.semanticscholar.org/CorpusID:258299717. (Cited on 7)

Xiangru Tang, Qiao Jin, Kunlun Zhu, Tongxin Yuan, Yichi Zhang, Wangchunshu Zhou, Meng Qu, Yilun Zhao, Jian Tang, Zhuosheng Zhang, Arman Cohan, Zhiyong Lu, and Mark B. Gerstein. Prioritizing Safeguarding Over Autonomy: Risks of LLM Agents for Science. *ArXiv*, abs/2402.04247, 2024. URL https://api.semanticscholar.org/CorpusID:267499916. (Cited on 22)

Simone Tedeschi, Felix Friedrich, Patrick Schramowski, Kristian Kersting, Roberto Navigli, Huu Nguyen, and Bo Li. ALERT: A Comprehensive Benchmark for Assessing Large Language Models' Safety through Red Teaming. *ArXiv*, abs/2404.08676, 2024. URL https://api.semanticscholar.org/CorpusID:269149567. (Cited on 3, 29)

Guardrails AI. Guardrails AI: Adding guardrails to large language models. https://github.com/guardrails-ai/guardrails, 2024. (Cited on 22, 27)

The White House. Executive Order on the Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence. 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/, 2023. (Cited on 2)

Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D. Manning, and Chelsea Finn. Fine-tuning Language Models for Factuality. *ArXiv*, abs/2311.08401, 2023. URL https://api.semanticscholar.org/CorpusID:265158181. (Cited on 24)

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. LLaMA: Open and Efficient Foundation Language Models. *ArXiv*, abs/2302.13971, 2023. URL https://api.semanticscholar.org/CorpusID:257219404. (Cited on 20)

Sam Toyer, Olivia Watkins, Ethan Adrian Mendes, Justin Svegliato, Luke Bailey, Tiffany Wang, Isaac Ong, Karim Elmaaroufi, Pieter Abbeel, Trevor Darrell, Alan Ritter, and Stuart Russell. Tensor Trust: Interpretable Prompt Injection Attacks from an Online Game. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=fsW7wJGLBd. (Cited on 2)

Florian Tramèr, Fan Zhang, Ari Juels, Michael K. Reiter, and Thomas Ristenpart. Stealing Machine Learning Models via Prediction APIs. In *USENIX Security Symposium*, 2016. URL https://api.semanticscholar.org/CorpusID:2984526. (Cited on 14)

Jean Marie Tshimula, Xavier Ndona, D'Jeff K. Nkashama, Pierre-Martin Tardif, Froduald Kabanza, Marc Frappier, and Shengrui Wang. Preventing jailbreak prompts as malicious tools for cybercriminals: A cyber defense perspective. *ArXiv*, abs/2411.16642, 2024. URL https://api.semanticscholar.org/CorpusID:274280511. (Cited on 3)

Alexander Matt Turner, Lisa Thiergart, David S. Udell, Gavin Leech, Ulisse Mini, and Monte Stuart MacDiarmid. Activation Addition: Steering Language Models Without Optimization. *ArXiv*, abs/2308.10248, 2023. URL https://api.semanticscholar.org/CorpusID:261049449. (Cited on 18)

Usage Policy Anthropic. Anthropic Usage Policies. https://www.anthropic.com/legal/aup, 2023. (Cited on 3)

Usage Policy OpenAI. OpenAI Usage Policies. https://openai.com/policies/usage-policies, 2023. (Cited on 3)

Thomas Vakili and Hercules Dalianis. Using Membership Inference Attacks to Evaluate Privacy-Preserving Language Modeling Fails for Pseudonymizing Data. In *Proceedings of the 24th Nordic Conference on Computational Linguistics (NoDaLiDa)*, pp. 318–323, 2023. (Cited on 15)

Matheus Valentim, Jeanette Falk, and Nanna Inie. Hacc-Man: An Arcade Game for Jailbreaking LLMs. 2024. URL https://api.semanticscholar.org/CorpusID:270063203. (Cited on 2)

Pranshu Verma and Will Oremus. These lawyers used ChatGPT to save time. They got fired and fined. https://www.washingtonpost.com/technology/2023/11/16/chatgpt-lawyer-fired-ai/, 2023. (Cited on 1)

Joe Vest and James Tubberville. Red Team Development and Operations–A practical Guide. *Independently Published*, 2020. (Cited on 1)

Bertie Vidgen, Adarsh Agrawal, Ahmed M. Ahmed, Victor Akinwande, Namir Al-nuaimi, Najla Alfaraj, Elie Alhajjar, Lora Aroyo, Trupti Bavalatti, Borhane Blili-Hamelin, Kurt D. Bollacker, Rishi Bomassani, Marisa Ferrara Boston, Sim'eon Campos, Kal Chakra, Canyu Chen, Cody Coleman, Zacharie Delpierre Coudert, Leon Derczynski, Debojyoti Dutta, Ian Eisenberg, James R. Ezick, Heather Frase, Brian Fuller, Ram Gandikota, Agasthya Gangavarapu, Ananya Gangavarapu, James Gealy, Rajat Ghosh, James Goel, Usman Gohar, Sujata Goswami, Scott A. Hale, Wiebke Hutiri, Joseph Marvin Imperial, Surgan Jandial, Nicholas C. Judd, Felix Juefei-Xu, Foutse Khomh, Bhavya Kailkhura, Hannah Rose Kirk, Kevin Klyman, Chris Knotz, Michael Kuchnik, Shachi H. Kumar, Chris Lengerich, Bo Li, Zeyi Liao, Eileen Peters Long, Victor Lu, Yifan Mai, Priyanka Mary Mammen, Kelvin Manyeki, Sean McGregor, Virendra

Mehta, Shafee Mohammed, Emanuel Moss, Lama Nachman, Dinesh Jinenhally Naganna, Amin Nikanjam, Besmira Nushi, Luis Oala, Iftach Orr, Alicia Parrish, Çigdem Patlak, William Pietri, Forough Poursabzi-Sangdeh, Eleonora Presani, Fabrizio Puletti, Paul Röttger, Saurav Sahay, Tim Santos, Nino Scherrer, Alice Schoenauer Sebag, Patrick Schramowski, Abolfazl Shahbazi, Vin Sharma, Xudong Shen, Vamsi Sistla, Leonard Tang, Davide Testuggine, Vithursan Thangarasa, Elizabeth Anne Watkins, Rebecca Weiss, Christoper A. Welty, Tyler Wilbers, Adina Williams, Carole-Jean Wu, Poonam Yadav, Xianjun Yang, Yi Zeng, Wenhui Zhang, Fedor Zhdanov, Jiacheng Zhu, Percy Liang, Peter Mattson, and Joaquin Vanschoren. Introducing v0.5 of the AI Safety Benchmark from MLCommons. *ArXiv*, abs/2404.12241, 2024. URL https://api.semanticscholar.org/CorpusID:269214329. (Cited on 3)

Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing NLP. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2153–2162, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1221. URL https://aclanthology.org/D19-1221. (Cited on 13)

Eric Wallace, Tony Zhao, Shi Feng, and Sameer Singh. Concealed Data Poisoning Attacks on NLP Models. In *North American Chapter of the Association for Computational Linguistics*, 2021. URL https://api.semanticscholar.org/CorpusID:233230124. (Cited on 19)

Eric Wallace, Kai Xiao, Reimar H. Leike, Lilian Weng, Johannes Heidecke, and Alex Beutel. The Instruction Hierarchy: Training LLMs to Prioritize Privileged Instructions. 2024. URL https://api.semanticscholar.org/CorpusID:269294048. (Cited on 24, 27)

Alexander Wan, Eric Wallace, Sheng Shen, and Dan Klein. Poisoning Language Models During Instruction Tuning. *ArXiv*, abs/2305.00944, 2023. URL https://api.semanticscholar.org/CorpusID:258426823. (Cited on 4, 19)

Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and B. Li. Adversarial GLUE: A Multi-Task Benchmark for Robustness Evaluation of Language Models. *ArXiv*, abs/2111.02840, 2021. URL https://api.semanticscholar.org/CorpusID:242757097. (Cited on 24)

Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, Sang Truong, Simran Arora, Mantas Mazeika, Dan Hendrycks, Zi-Han Lin, Yuk-Kit Cheng, Sanmi Koyejo, Dawn Xiaodong Song, and Bo Li. DecodingTrust: A Comprehensive Assessment of Trustworthiness in GPT Models. *ArXiv*, abs/2306.11698, 2023a. URL https://api.semanticscholar.org/CorpusID:259202782. (Cited on 3, 5, 6, 29)

Chaoqi Wang, Yibo Jiang, Chenghao Yang, Han Liu, and Yuxin Chen. Beyond reverse kl: Generalizing direct preference optimization with diverse divergence constraints. *ArXiv*, abs/2309.16240, 2023b. URL https://api.semanticscholar.org/CorpusID:263142109. (Cited on 24)

Haoran Wang and Kai Shu. Backdoor Activation Attack: Attack Large Language Models using Activation Steering for Safety-Alignment. *ArXiv*, abs/2311.09433, 2023. URL https://api.semanticscholar.org/CorpusID:265220823. (Cited on 4, 18)

Jiong Wang, Zi yang Liu, Keun Hee Park, Muhao Chen, and Chaowei Xiao. Adversarial Demonstration Attacks on Large Language Models. *ArXiv*, abs/2305.14950, 2023c. URL https://api.semanticscholar.org/CorpusID:258865399. (Cited on 4)

Jiong Wang, Jiazhao Li, Yiquan Li, Xiangyu Qi, Junjie Hu, Yixuan Li, Patrick Drew McDaniel, Muhao Chen, Bo Li, and Chaowei Xiao. Mitigating Fine-tuning Jailbreak Attack with Backdoor Enhanced Alignment. *ArXiv*, abs/2402.14968, 2024a. URL https://api.semanticscholar.org/CorpusID:267897454. (Cited on 25, 27)

Lei Wang, Chengbang Ma, Xueyang Feng, Zeyu Zhang, Hao ran Yang, Jingsen Zhang, Zhi-Yang Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji rong Wen. A Survey on Large Language

Model based Autonomous Agents. *ArXiv*, abs/2308.11432, 2023d. URL https://api.semanticscholar.org/CorpusID:261064713. (Cited on 7, 22)

Meng Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi Yang, Jindong Wang, and Huajun Chen. Detoxifying Large Language Models via Knowledge Editing. *ArXiv*, abs/2403.14472, 2024b. URL https://api.semanticscholar.org/CorpusID:268553537. (Cited on 28, 29)

Yihan Wang, Zhouxing Shi, Andrew Bai, and Cho-Jui Hsieh. Defending LLMs against Jailbreaking Attacks via Backtranslation. *ArXiv*, abs/2402.16459, 2024c. URL https://api.semanticscholar.org/CorpusID:268032484. (Cited on 22, 27)

Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. Aligning Large Language Models with Human: A Survey. *ArXiv*, abs/2307.12966, 2023e. URL https://api.semanticscholar.org/CorpusID:260356605. (Cited on 2)

Ziqiu Wang, Jun Liu, Shengkai Zhang, and Yang Yang. Poisoned LangChain: Jailbreak LLMs by LangChain. 2024d. URL https://api.semanticscholar.org/CorpusID:270738192. (Cited on 18)

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How Does LLM Safety Training Fail? *ArXiv*, abs/2307.02483, 2023a. URL https://api.semanticscholar.org/CorpusID:259342528. (Cited on 4, 11)

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent Abilities of Large Language Models. *Trans. Mach. Learn. Res.*, 2022, 2022. URL https://openreview.net/forum?id=yzkSU5zdwD. (Cited on 1)

Junyin Wei, Yicheng Zhang, Zhe Zhou, Zhou Li, and Mohammad Abdullah Al Faruque. Leaky DNN: Stealing Deep-Learning Model Secret with GPU Context-Switching Side-Channel. *2020 50th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, pp. 125–137, 2020. URL https://api.semanticscholar.org/CorpusID:220939286. (Cited on 17)

Lingxiao Wei, Yannan Liu, Bo Luo, Yu LI, and Qiang Xu. I Know What You See: Power Side-Channel Attack on Convolutional Neural Network Accelerators. *Proceedings of the 34th Annual Computer Security Applications Conference*, 2018. URL https://api.semanticscholar.org/CorpusID:3920864. (Cited on 17)

Zeming Wei, Yifei Wang, and Yisen Wang. Jailbreak and Guard Aligned Language Models with Only Few In-Context Demonstrations. *ArXiv*, abs/2310.06387, 2023b. URL https://api.semanticscholar.org/CorpusID:263830179. (Cited on 4)

Laura Weidinger, John F. J. Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zachary Kenton, Sande Minnich Brown, William T. Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William S. Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Ethical and social risks of harm from Language Models. *ArXiv*, abs/2112.04359, 2021. URL https://api.semanticscholar.org/CorpusID:244954639. (Cited on 1, 6, 29)

Laura Weidinger, John F J Mellor, Bernat Guillén Pegueroles, Nahema Marchal, Ravin Kumar, Kristian Lum, Canfer Akbulut, Mark Diaz, A. Stevie Bergman, Mikel D. Rodriguez, Verena Rieser, and William Isaac. STAR: SocioTechnical approach to red teaming language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 21516–21532, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1200. URL https://aclanthology.org/2024.emnlp-main.1200/. (Cited on 31)

Roy Weiss, Daniel Ayzenshteyn, Guy Amit, and Yisroel Mirsky. What Was Your Prompt? A Remote Keylogging Attack on AI Assistants. 2024. URL https://api.semanticscholar.org/CorpusID:268510380. (Cited on 4)

Jiaxin Wen, Pei Ke, Hao Sun, Zhexin Zhang, Chengfei Li, Jinfeng Bai, and Minlie Huang. Unveiling the Implicit Toxicity in Large Language Models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1322–1338, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.84. URL https://aclanthology.org/2023.emnlp-main.84. (Cited on 4)

Nevan Wichers, Carson E. Denison, and Ahmad Beirami. Gradient-Based Language Model Red Teaming. *ArXiv*, abs/2401.16656, 2024. URL https://api.semanticscholar.org/CorpusID:267320331. (Cited on 4, 21)

David Gray Widder, Sarah West, and Meredith Whittaker. Open (For Business): Big Tech, Concentrated Power, and the Political Economy of Open AI. *SSRN Electronic Journal*, 2023. URL https://api.semanticscholar.org/CorpusID:265667031. (Cited on 2)

Yotam Wolf, Noam Wies, Yoav Levine, and Amnon Shashua. Fundamental Limitations of Alignment in Large Language Models. *ArXiv*, abs/2304.11082, 2023. URL https://api.semanticscholar.org/CorpusID:258291526. (Cited on 30)

Baoyuan Wu, Hongrui Chen, Mingda Zhang, Zihao Zhu, Shaokui Wei, Danni Yuan, Mingli Zhu, Ruotong Wang, Li Liu, and Chaoxiao Shen. BackdoorBench: A Comprehensive Benchmark and Analysis of Backdoor Learning. *ArXiv*, abs/2401.15002, 2024a. URL https://api.semanticscholar.org/CorpusID:267301142. (Cited on 29)

Xi Wu, Matthew Fredrikson, Somesh Jha, and Jeffrey F Naughton. A methodology for formalizing model-inversion attacks. In *2016 IEEE 29th Computer Security Foundations Symposium (CSF)*, pp. 355–370. IEEE, 2016. (Cited on 4, 15)

Yuanwei Wu, Xiang Li, Yixin Liu, Pan Zhou, and Lichao Sun. Jailbreaking GPT-4V via Self-Adversarial Attacks with System Prompts. *CoRR*, abs/2311.09127, 2023. doi: 10.48550/ARXIV.2311.09127. URL https://doi.org/10.48550/arXiv.2311.09127. (Cited on 4)

Yuhao Wu, Franziska Roesner, Tadayoshi Kohno, Ning Zhang, and Umar Iqbal. SecGPT: An Execution Isolation Architecture for LLM-Based Systems. 2024b. URL https://api.semanticscholar.org/CorpusID:268296997. (Cited on 31)

Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. OS-Copilot: Towards Generalist Computer Agents with Self-Improvement. *ArXiv*, abs/2402.07456, 2024c. URL https://api.semanticscholar.org/CorpusID:267626905. (Cited on 7)

Chong Xiang, Tong Wu, Zexuan Zhong, David A. Wagner, Danqi Chen, and Prateek Mittal. Certifiably robust RAG against retrieval corruption. *CoRR*, abs/2405.15556, 2024a. doi: 10.48550/ARXIV.2405.15556. URL https://doi.org/10.48550/arXiv.2405.15556. (Cited on 18)

Yun Xiang, Zhuangzhi Chen, Zuohui Chen, Zebin Fang, Haiyang Hao, Jinyin Chen, Yi Liu, Zhefu Wu, Qi Xuan, and Xiaoniu Yang. Open DNN Box by Power Side-Channel Attack. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 67:2717–2721, 2019. URL https://api.semanticscholar.org/CorpusID:198229526. (Cited on 17)

Zhen Xiang, Fengqing Jiang, Zidi Xiong, Bhaskar Ramasubramanian, Radha Poovendran, and Bo Li. BadChain: Backdoor Chain-of-Thought Prompting for Large Language Models. *ArXiv*, abs/2401.12242, 2024b. URL https://api.semanticscholar.org/CorpusID:267094837. (Cited on 4, 18)

Yijia Xiao, Yiqiao Jin, Yushi Bai, Yue Wu, Xianjun Yang, Xiao Luo, Wenchao Yu, Xujiang Zhao, Yanchi Liu, Haifeng Chen, Wei Wang, and Wei Cheng. Large Language Models Can Be Good Privacy Protection Learners. *ArXiv*, abs/2310.02469, 2023. URL https://api.semanticscholar.org/CorpusID:263620236. (Cited on 27)

Zeguan Xiao, Yan Yang, Guanhua Chen, and Yun Chen. Tastle: Distract Large Language Models for Automatic Jailbreak Attack. *ArXiv*, abs/2403.08424, 2024. URL https://api.semanticscholar.org/CorpusID:268379559. (Cited on 28)

Tinghao Xie, Xiangyu Qi, Yi Zeng, Yangsibo Huang, Udari Madhushani Sehwag, Kaixuan Huang, Luxi He, Boyi Wei, Dacheng Li, Ying Sheng, Ruoxi Jia, Bo Li, Kai Li, Danqi Chen, Peter Henderson, and Prateek Mittal. SORRY-Bench: Systematically Evaluating Large Language Model Safety Refusal Behaviors. 2024a. URL https://api.semanticscholar.org/CorpusID:270688409. (Cited on 29, 31)

Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl, Lingjuan Lyu, Qifeng Chen, Xing Xie, and Fangzhao Wu. Defending ChatGPT against jailbreak attack via self-reminders. *Nature Machine Intelligence*, 5:1486–1496, 2023. URL https://api.semanticscholar.org/CorpusID:266289038. (Cited on 22, 27)

Yueqi Xie, Minghong Fang, Renjie Pi, and Neil Zhenqiang Gong. GradSafe: Detecting Jailbreak Prompts for LLMs via Safety-Critical Gradient Analysis. 2024b. URL https://api.semanticscholar.org/CorpusID:267770418. (Cited on 28)

Guo Xingang, Fangxu Yu, Huan Zhang, Lianhui Qin, and Bin Hu. COLD-Attack: Jailbreaking LLMs with Stealthiness and Controllability. *ArXiv*, abs/2402.08679, 2024. URL https://api.semanticscholar.org/CorpusID:267637251. (Cited on 4, 21)

Chen Xiong, Xiangyu Qi, Pin-Yu Chen, and Tsung-Yi Ho. Defensive Prompt Patch: A Robust and Interpretable Defense of LLMs against Jailbreak Attacks. 2024. URL https://api.semanticscholar.org/CorpusID:270123437. (Cited on 27)

Jiashu Xu, Mingyu Derek Ma, Fei Wang, Chaowei Xiao, and Muhao Chen. Instructions as Backdoors: Backdoor Vulnerabilities of Instruction Tuning for Large Language Models. *ArXiv*, abs/2305.14710, 2023a. URL https://api.semanticscholar.org/CorpusID:258866212. (Cited on 4, 19)

Nan Xu, Fei Wang, Ben Zhou, Bangzheng Li, Chaowei Xiao, and Muhao Chen. Cognitive Overload: Jailbreaking Large Language Models with Overloaded Logical Thinking. *ArXiv*, abs/2311.09827, 2023b. URL https://api.semanticscholar.org/CorpusID:265221395. (Cited on 4, 13)

Yue Xu and Wenjie Wang. LinkPrompt: Natural and Universal Adversarial Attacks on Prompt-based Language Models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 6473–6486, Mexico City, Mexico, June 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.naacl-long.360. (Cited on 14)

Zhangchen Xu, Fengqing Jiang, Luyao Niu, Jinyuan Jia, Bill Yuchen Lin, and Radha Poovendran. SafeDecoding: Defending against Jailbreak Attacks via Safety-Aware Decoding. *ArXiv*, abs/2402.08983, 2024a. URL https://api.semanticscholar.org/CorpusID:267658033. (Cited on 23, 27)

Zihao Xu, Yi Liu, Gelei Deng, Yuekang Li, and Stjepan Picek. LLM Jailbreak Attack versus Defense Techniques – A Comprehensive Study. 2024b. URL https://api.semanticscholar.org/CorpusID:267770234. (Cited on 22)

Jiaqi Xue, Mengxin Zheng, Ting Hua, Yilin Shen, Yepeng Liu, Ladislau Bölöni, and Qian Lou. TrojLLM: A Black-box Trojan Prompt Attack on Large Language Models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 65665–65677. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/cf04d01a0e76f8b13095349d9caca033-Paper-Conference.pdf. (Cited on 4, 13)

Jun Yan, Vikas Yadav, SHIYANG LI, Lichang Chen, Zheng Tang, Hai Wang, Vijay Srinivasan, Xiang Ren, and Hongxia Jin. Backdooring Instruction-Tuned Large Language Models with Virtual Prompt Injection. 2023. URL https://api.semanticscholar.org/CorpusID:260334112. (Cited on 4)

Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond. *ArXiv*, abs/2304.13712, 2023a. URL https://api.semanticscholar.org/CorpusID:258331833. (Cited on 3)

Wenkai Yang, Lei Li, Zhiyuan Zhang, Xuancheng Ren, Xu Sun, and Bin He. Be Careful about Poisoned Word Embeddings: Exploring the Vulnerability of the Embedding Layers in NLP Models. *ArXiv*, abs/2103.15543, 2021a. URL https://api.semanticscholar.org/CorpusID:232404131. (Cited on 19)

Wenkai Yang, Yankai Lin, Peng Li, Jie Zhou, and Xu Sun. Rethinking Stealthiness of Backdoor Attack against NLP Models. In *Annual Meeting of the Association for Computational Linguistics*, 2021b. URL https://api.semanticscholar.org/CorpusID:236459933. (Cited on 19)

Wenkai Yang, Xiaohan Bi, Yankai Lin, Sishuo Chen, Jie Zhou, and Xu Sun. Watch Out for Your Agents! Investigating Backdoor Threats to LLM-Based Agents. *ArXiv*, abs/2402.11208, 2024. URL https://api.semanticscholar.org/CorpusID:267751034. (Cited on 22)

Xianjun Yang, Xiao Wang, Qi Zhang, Linda Ruth Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. Shadow Alignment: The Ease of Subverting Safely-Aligned Language Models. *ArXiv*, abs/2310.02949, 2023b. URL https://api.semanticscholar.org/CorpusID:263620436. (Cited on 4, 20)

Yuchen Yang, Bo Hui, Haolin Yuan, Neil Gong, and Yinzhi Cao. SneakyPrompt: Jailbreaking Text-to-image Generative Models. 2023c. URL https://api.semanticscholar.org/CorpusID:265150147. (Cited on 6)

Zhilin Yang, Zihang Dai, Ruslan Salakhutdinov, and William W. Cohen. Breaking the Softmax Bottleneck: A High-Rank RNN Language Model. *ArXiv*, abs/1711.03953, 2017. URL https://api.semanticscholar.org/CorpusID:26238954. (Cited on 15)

Zhou Yang, Bowen Xu, J Zhang, Hong Jin Kang, Jieke Shi, Junda He, and David Lo. Stealthy Backdoor Attack for Code Models. *ArXiv*, abs/2301.02496, 2023d. URL https://api.semanticscholar.org/CorpusID:255522586. (Cited on 19)

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing Reasoning and Acting in Language Models. *ArXiv*, abs/2210.03629, 2022. URL https://api.semanticscholar.org/CorpusID:252762395. (Cited on 7)

Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large Language Model Unlearning. *CoRR*, abs/2310.10683, 2023. doi: 10.48550/ARXIV.2310.10683. URL https://doi.org/10.48550/arXiv.2310.10683. (Cited on 4, 20)

Dayong Ye, Sheng Shen, Tianqing Zhu, Bo Liu, and Wanlei Zhou. One parameter defense—defending against data inference attacks via differential privacy. *IEEE Transactions on Information Forensics and Security*, 17:1466–1480, 2022. (Cited on 25)

Mao Ye, Chengyue Gong, and Qiang Liu. SAFER: A Structure-free Approach for Certified Robustness to Adversarial Word Substitutions. In *Annual Meeting of the Association for Computational Linguistics*, 2020. URL https://api.semanticscholar.org/CorpusID:219124328. (Cited on 23)

Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy Risk in Machine Learning: Analyzing the Connection to Overfitting. *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, pp. 268–282, 2017. URL https://api.semanticscholar.org/CorpusID:2656445. (Cited on 14)

Jingwei Yi, Yueqi Xie, Bin Zhu, Keegan Hines, Emre Kiciman, Guangzhong Sun, Xing Xie, and Fangzhao Wu. Benchmarking and Defending Against Indirect Prompt Injection Attacks on Large Language Models. *ArXiv*, abs/2312.14197, 2023. URL https://api.semanticscholar.org/CorpusID:266521508. (Cited on 29)

Xin Yi, Shunfan Zheng, Linlin Wang, Xiaoling Wang, and Liang He. A safety realignment framework via subspace-oriented model fusion for large language models. *ArXiv*, abs/2405.09055, 2024. URL https://api.semanticscholar.org/CorpusID:269773206. (Cited on 27)

Itay Yona, Ilia Shumailov, Jamie Hayes, and Nicholas Carlini. Stealing user prompts from mixture of experts. *CoRR*, abs/2410.22884, 2024. doi: 10.48550/ARXIV.2410.22884. URL https://doi.org/10.48550/arXiv.2410.22884. (Cited on 17)

Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A. Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, Sergey Yekhanin, and Huishuai Zhang. Differentially Private Fine-tuning of Language Models. *ArXiv*, abs/2110.06500, 2021a. URL https://api.semanticscholar.org/CorpusID:238743879. (Cited on 25, 27)

Da Yu, Huishuai Zhang, Wei Chen, Jian Yin, and Tie-Yan Liu. Large Scale Private Learning via Low-rank Reparametrization. *ArXiv*, abs/2106.09352, 2021b. URL https://api.semanticscholar.org/CorpusID:235458199. (Cited on 25)

Jiahao Yu, Xingwei Lin, Zheng Yu, and Xinyu Xing. GPTFUZZER: Red Teaming Large Language Models with Auto-Generated Jailbreak Prompts. *ArXiv*, abs/2309.10253, 2023a. URL https://api.semanticscholar.org/CorpusID:262055242. (Cited on 4, 13)

Weichen Yu, Tianyu Pang, Qian Liu, Chao Du, Bingyi Kang, Yan Huang, Min Lin, and Shuicheng Yan. Bag of Tricks for Training Data Extraction from Language Models. *ArXiv*, abs/2302.04460, 2023b. URL https://api.semanticscholar.org/CorpusID:256697118. (Cited on 4, 15)

Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen tse Huang, Pinjia He, Shuming Shi, and Zhaopeng Tu. GPT-4 Is Too Smart To Be Safe: Stealthy Chat with LLMs via Cipher. *ArXiv*, abs/2308.06463, 2023. URL https://api.semanticscholar.org/CorpusID:260887189. (Cited on 4, 12)

Zhuowen Yuan, Zidi Xiong, Yi Zeng, Ning Yu, Ruoxi Jia, Dawn Xiaodong Song, and Bo Li. RigorLLM: Resilient Guardrails for Large Language Models against Undesired Content. 2024. URL https://api.semanticscholar.org/CorpusID:268536710. (Cited on 22, 27)

Canaan Yung, Hadi Mohaghegh Dolatabadi, Sarah Monazam Erfani, and Christopher Leckie. Round Trip Translation Defence against Large Language Model Jailbreaking Attacks. 2024. URL https://api.semanticscholar.org/CorpusID:267770468. (Cited on 22, 27)

Matei Zaharia, Omar Khattab, Lingjiao Chen, Jared Quincy Davis, Heather Miller, Chris Potts, James Zou, Michael Carbin, Jonathan Frankle, Naveen Rao, and Ali Ghodsi. The Shift from Models to Compound AI Systems. https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/, 2024. (Cited on 7, 22)

Yi Zeng, Kevin Klyman, Andy Zhou, Yu Yang, Minzhou Pan, Ruoxi Jia, Dawn Song, Percy Liang, and Bo Li. Ai risk categorization decoded (air 2024): From government regulations to corporate policies. *ArXiv*, abs/2406.17864, 2024a. URL https://api.semanticscholar.org/CorpusID:270738014. (Cited on 3)

Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. How Johnny Can Persuade LLMs to Jailbreak Them: Rethinking Persuasion to Challenge AI Safety by Humanizing LLMs. *ArXiv*, abs/2401.06373, 2024b. URL https://api.semanticscholar.org/CorpusID:266977395. (Cited on 4)

Yi Zeng, Weiyu Sun, Tran Ngoc Huynh, Dawn Song, Bo Li, and Ruoxi Jia. BEEAR: Embedding-based Adversarial Removal of Safety Backdoors in Instruction-tuned Language Models. 2024c. URL https://api.semanticscholar.org/CorpusID:270711308. (Cited on 25)

Yifan Zeng, Yiran Wu, Xiao Zhang, Huazheng Wang, and Qingyun Wu. AutoDefense: Multi-Agent LLM Defense against Jailbreak Attacks. *ArXiv*, abs/2403.04783, 2024d. URL https://api.semanticscholar.org/CorpusID:268297202. (Cited on 26)

Qiusi Zhan, Richard Fang, R. Tanya Bindu, Akul Gupta, Tatsunori Hashimoto, and Daniel Kang. Removing RLHF Protections in GPT-4 via Fine-Tuning. *ArXiv*, abs/2311.05553, 2023. URL https://api.semanticscholar.org/CorpusID:265067269. (Cited on 4, 20)

Jie Zhang, Haoyu Bu, Hui Wen, Yu Chen, Lun Li, and Hongsong Zhu. When LLMs Meet Cybersecurity: A Systematic Literature Review. 2024a. URL https://api.semanticscholar.org/CorpusID:269604648. (Cited on 6)

Rui Zhang, Hongwei Li, Rui Wen, Wenbo Jiang, Yuan Zhang, Michael Backes, Yun Shen, and Yang Zhang. Rapid Adoption, Hidden Risks: The Dual Impact of Large Language Model Customization. 2024b. URL https://api.semanticscholar.org/CorpusID:267657775. (Cited on 2)

Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang, Kathleen McKeown, and Tatsunori Hashimoto. Benchmarking Large Language Models for News Summarization. *Transactions of the Association for Computational Linguistics*, 12:39–57, 2023a. URL https://api.semanticscholar.org/CorpusID:256416014. (Cited on 7)

Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. Self-Alignment for Factuality: Mitigating Hallucinations in LLMs via Self-Evaluation. *ArXiv*, abs/2402.09267, 2024c. URL https://api.semanticscholar.org/CorpusID:267657805. (Cited on 24)

Xinyang Zhang, Zheng Zhang, and Ting Wang. Trojaning Language Models for Fun and Profit. *2021 IEEE European Symposium on Security and Privacy (EuroS&P)*, pp. 179–197, 2020. URL https://api.semanticscholar.org/CorpusID:220936152. (Cited on 19)

Yihao Zhang and Zeming Wei. Boosting Jailbreak Attack with Momentum. 2024. URL https://api.semanticscholar.org/CorpusID:269502544. (Cited on 4)

Yiming Zhang and Daphne Ippolito. Effective Prompt Extraction from Language Models. 2023. URL https://api.semanticscholar.org/CorpusID:259847681. (Cited on 4)

Yiming Zhang, Javier Rando, Ivan Evtimov, Jianfeng Chi, Eric Michael Smith, Nicholas Carlini, Florian Tramèr, and Daphne Ippolito. Persistent pre-training poisoning of llms. *CoRR*, abs/2410.13722, 2024d. doi: 10.48550/ARXIV.2410.13722. URL https://doi.org/10.48550/arXiv.2410.13722. (Cited on 4, 19)

Yuqi Zhang, Liang Ding, Lefei Zhang, and Dacheng Tao. Intention Analysis Prompting Makes Large Language Models A Good Jailbreak Defender. 2024e. URL https://api.semanticscholar.org/CorpusID:266977251. (Cited on 22, 27)

Zhen Zhang, Guanhua Zhang, Bairu Hou, Wenqi Fan, Qing Li, Sijia Liu, Yang Zhang, and Shiyu Chang. Certified Robustness for Large Language Models with Self-Denoising. *ArXiv*, abs/2307.07171, 2023b. URL https://api.semanticscholar.org/CorpusID:259924952. (Cited on 23)

Zhengyan Zhang, Guangxuan Xiao, Yongwei Li, Tian Lv, Fanchao Qi, Yasheng Wang, Xin Jiang, Zhiyuan Liu, and Maosong Sun. Red Alarm for Pre-trained Models: Universal Vulnerability to Neuron-level Backdoor Attacks. *Machine Intelligence Research*, 20:180–193, 2021. URL https://api.semanticscholar.org/CorpusID:231632260. (Cited on 19)

Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. SafetyBench: Evaluating the Safety of Large Language Models with Multiple Choice Questions. *ArXiv*, abs/2309.07045, 2023c. URL https://api.semanticscholar.org/CorpusID:261706197. (Cited on 29)

Zhexin Zhang, Junxiao Yang, Pei Ke, and Minlie Huang. Defending Large Language Models Against Jailbreaking Attacks Through Goal Prioritization. *ArXiv*, abs/2311.09096, 2023d. URL https://api.semanticscholar.org/CorpusID:265212812. (Cited on 24, 27)

Zhuo Zhang, Guangyu Shen, Guanhong Tao, Siyuan Cheng, and Xiangyu Zhang. Make Them Spill the Beans! Coercive Knowledge Extraction from (Production) LLMs. *ArXiv*, abs/2312.04782, 2023e. URL https://api.semanticscholar.org/CorpusID:266149700. (Cited on 4, 13)

Haiteng Zhao, Chang Ma, Xinshuai Dong, Anh Tuan Luu, Zhi-Hong Deng, and Hanwang Zhang. Certified robustness against natural language attacks by causal intervention. In *International Conference on Machine Learning*, pp. 26958–26970. PMLR, 2022. (Cited on 23)

Shuai Zhao, Meihuizi Jia, Anh Tuan Luu, and Jinming Wen. Universal Vulnerabilities in Large Language Models: In-context Learning Backdoor Attacks. *ArXiv*, abs/2401.05949, 2024a. URL https://api.semanticscholar.org/CorpusID:266932984. (Cited on 4, 18)

Wei Zhao, Zhe Li, Yige Li, Ye Zhang, and Junfeng Sun. Defending Large Language Models Against Jailbreak Attacks via Layer-specific Editing. 2024b. URL https://api.semanticscholar.org/CorpusID:270067915. (Cited on 27)

Xuandong Zhao, Xianjun Yang, Tianyu Pang, Chao Du, Lei Li, Yu-Xiang Wang, and William Yang Wang. Weak-to-Strong Jailbreaking on Large Language Models. *ArXiv*, abs/2401.17256, 2024c. URL https://api.semanticscholar.org/CorpusID:267320277. (Cited on 4)

Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J. Liu. Slic-hf: Sequence likelihood calibration with human feedback. *ArXiv*, abs/2305.10425, 2023. URL https://api.semanticscholar.org/CorpusID:258741082. (Cited on 24)

Andy Zhou, Bo Li, and Haohan Wang. Robust Prompt Optimization for Defending Language Models Against Jailbreaking Attacks. *ArXiv*, abs/2401.17263, 2024a. URL https://api.semanticscholar.org/CorpusID:267320750. (Cited on 22, 27)

Jingyan Zhou, Kun Li, Junan Li, Jiawen Kang, Minda Hu, Xixin Wu, and Helen M. Meng. Purple-teaming LLMs with Adversarial Defender Training. 2024b. URL https://api.semanticscholar.org/CorpusID:270878226. (Cited on 28)

YongBin Zhou and DengGuo Feng. Side-channel attacks: Ten years after its publication and the impacts on cryptographic module security testing. *Cryptology ePrint Archive*, 2005. (Cited on 16)

Banghua Zhu, Norman Mu, Jiantao Jiao, and David Wagner. Generative AI Security: Challenges and Countermeasures. 2024. URL https://api.semanticscholar.org/CorpusID:267759992. (Cited on 29)

Sicheng Zhu, Ruiyi Zhang, Bang An, Gang Wu, Joe Barrow, Zichao Wang, Furong Huang, Ani Nenkova, and Tong Sun. AutoDAN: Interpretable Gradient-Based Adversarial Attacks on Large Language Models. 2023. URL https://api.semanticscholar.org/CorpusID:264451545. (Cited on 4, 10, 14, 21)

Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. Red teaming ChatGPT via Jailbreaking: Bias, Robustness, Reliability and Toxicity. 2023. URL https://api.semanticscholar.org/CorpusID:258960373. (Cited on 6)

Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-Tuning Language Models from Human Preferences. *arXiv preprint arXiv:1909.08593*, 2019. URL https://arxiv.org/abs/1909.08593. (Cited on 24)

Daniel M. Ziegler, Seraphina Nix, Lawrence Chan, Tim Bauman, Peter Schmidt-Nielsen, Tao Lin, Adam Scherlis, Noa Nabeshima, Ben Weinstein-Raun, Daniel Haas, Buck Shlegeris, and Nate Thomas. Adversarial Training for High-Stakes Reliability. *ArXiv*, abs/2205.01663, 2022. URL https://api.semanticscholar.org/CorpusID:248506146. (Cited on 24)

Caleb Ziems, William B. Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. Can Large Language Models Transform Computational Social Science? *ArXiv*, abs/2305.03514, 2023. URL https://api.semanticscholar.org/CorpusID:258547324. (Cited on 7)

Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and Transferable Adversarial Attacks on Aligned Language Models. *ArXiv*, abs/2307.15043, 2023. URL https://api.semanticscholar.org/CorpusID:260202961. (Cited on 4, 10, 13, 14, 21, 29)

Wei Zou, Runpeng Geng, Binghui Wang, and Jinyuan Jia. PoisonedRAG: Knowledge Poisoning Attacks to Retrieval-Augmented Generation of Large Language Models. *ArXiv*, abs/2402.07867, 2024. URL https://api.semanticscholar.org/CorpusID:267626957. (Cited on 4, 18, 23, 28)