



BRNO UNIVERSITY OF TECHNOLOGY

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

FACULTY OF INFORMATION TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

DEPARTMENT OF COMPUTER GRAPHICS AND MULTIMEDIA

ÚSTAV POČÍTAČOVÉ GRAFIKY A MULTIMÉDIÍ

DEVELOPMENT OF AN LLM RED-TEAMING TOOLKIT

VÝVOJ TOOLKITU PRO RED-TEAMING VELKÝCH JAZYKOVÝCH MODELŮ (LLM)

BACHELOR'S THESIS

BAKALÁŘSKÁ PRÁCE

AUTHOR

AUTOR PRÁCE

SUPERVISOR

VEDOUCÍ PRÁCE

ADAM VESELÝ

Ing. JAKUB REŠ

BRNO 2026

Abstract

Do tohoto odstavce bude zapsán výtah (abstrakt) práce v anglickém jazyce.

Abstrakt

Do tohoto odstavce bude zapsán výtah (abstrakt) práce v českém (slovenském) jazyce.

Keywords

Sem budou zapsána jednotlivá klíčová slova v anglickém jazyce, oddělená čárkami.

Klíčová slova

Sem budou zapsána jednotlivá klíčová slova v českém (slovenském) jazyce, oddělená čárkami.

Reference

VESELÝ, Adam. *Development of an LLM Red-Teaming Toolkit*. Brno, 2026. Bachelor's thesis. Brno University of Technology, Faculty of Information Technology. Supervisor Ing. Jakub Reš,

Development of an LLM Red-Teaming Toolkit

Declaration

I hereby declare that this Bachelor's thesis was prepared as an original work by the author under the supervision of Mr. Ing. Jakub Reš. I have listed all the literary sources, publications and other sources, which were used during the preparation of this thesis. I have used ChatGPT to correct spelling and other language mistakes. I used Grok when working on the software.

.....
Adam Veselý
January 6, 2026

Acknowledgements

I would like to thank my supervisor Ing. Jakub Reš for his guidance and support throughout the development of this thesis. I also appreciate the assistance provided by my colleagues and friends who contributed their insights and expertise.

Contents

1	Introduction	3
2	Background on Large Language Models and Safety Risks	4
2.1	Large Language Models	4
2.2	LLM Safety Risks	5
3	Red-Teaming of Large Language Models	7
3.1	Red-teaming: Definitions and Methodology	7
3.2	Attack Taxonomy	8
3.3	Evaluation Metrics and Benchmarks	13
3.4	Open-source Red-Teaming Tools and Frameworks	14
3.5	Synthesis and Research Gap	17
4	Design	18
4.1	Requirements and Goals	18
4.2	High-level Architecture	19
4.3	Module Specifications	19
4.4	Default Test Suite and OWASP Mapping	19
4.5	Security, Privacy and Ethical Considerations	19
4.6	Proof-of-Concept Deployment	19
5	Implementation	20
6	Evaluation	21
7	Conclusion	22
	Bibliography	23

List of Figures

3.1	Example of a single-turn jailbreak prompt that elicits disallowed content from an LLM in one prompt. Adapted from Zhao and Zhang [32].	9
3.2	Example of an indirect prompt injection attack: a user issues a benign query, the LLM retrieves external content containing embedded instructions, and these instructions override system or developer prompts, leading to compromised output. Adapted from Greshake et al. [7].	9
3.3	Illustration of a multi-turn jailbreak attack, in which an attacker iteratively refines prompts based on previous model responses to progressively elicit disallowed behavior. Adapted from Zhao and Zhang [32].	10
3.4	Illustration of a universal trigger attack, where a fixed trigger sequence appended to user inputs induces undesired behavior in an unmodified language model across different prompts and, in some cases, across different models. .	11
3.5	Illustration of a Trojan-style (backdoor) attack on a large language model. The modified model behaves normally for benign inputs but produces malicious or policy-violating outputs when a specific trigger pattern is present. Unlike universal trigger attacks, the vulnerability is embedded in the model during training or fine-tuning rather than in the input prompt.	12
3.6	Illustration of a prompt-based evasion attack, where an attacker modifies a disallowed request using linguistic obfuscation or adversarial perturbations to bypass safety mechanisms.	13
4.1	High-level architecture of the proposed red-teaming toolkit. The system is organised around a central orchestrator that coordinates attack generators, model adapters, and a hybrid judging pipeline. Experiments are defined declaratively via manifests and logged in a reproducible JSONL format, with optional indexing and reporting via a lightweight UI.	18

Chapter 1

Introduction

The rapid advancement and widespread deployment of large language models (LLMs) have transformed natural-language interaction with computers. These models now power chatbots, code assistants, translation systems, and creative tools used daily by millions of users.

However, their remarkable capabilities come with significant safety and ethical risks. LLMs can generate harmful, biased, misleading, or illegal content when subjected to carefully crafted adversarial inputs, a practice commonly known as *jailbreaking* [29, 33].

Real-world incidents such as ChatGPT being tricked into providing bomb-making instructions [5], or Gemini’s image-generation controversy [26], have demonstrated that even flagship commercial models remain vulnerable to adversarial prompting. In response, red-teaming, a cybersecurity technique involving simulated attacks to expose vulnerabilities, has been adopted by leading AI organisations (OpenAI, Anthropic, Google DeepMind) as a core component of LLM safety evaluation [16].

Beyond model-level behaviour, LLMs are embedded in applications and services where integration and deployment issues create additional, system-level risks. Industry guidance such as the OWASP Top 10 for Large Language Model Applications highlights vulnerabilities including prompt injection, insecure output handling, insecure plugin design, and excessive agency [18, 19].

With the adoption of the EU AI Act in 2024, systematic risk assessment including red-teaming will become a legal requirement for high-risk AI systems deployed in the European Union from 2026 onward [20]. Consequently, efficient, reproducible, and extensible red-teaming tools are no longer a luxury but an essential part of responsible AI development.

Despite significant progress, most existing open-source red-teaming frameworks suffer from limited modularity, poor support for modern systems and computational requirements that hinder adoption by smaller research teams and individual developers [22, 24, 1, 15]. This creates a clear need for a new, lightweight, developer-friendly red-teaming toolkit that lowers the barrier to LLM systematic safety testing.

The main goal of this bachelor’s thesis is therefore the design, implementation, and evaluation of a modular open-source red-teaming toolkit for large language models that addresses the identified shortcomings of current solutions.

The remainder of this thesis is structured as follows: Chapter 2 provides the necessary background on large language models. Chapter 3 surveys existing red-teaming methodologies, attack taxonomies, evaluation metrics, and open-source tools. Chapter 4 presents the proposed system architecture. Chapter 5 details the implementation of the red-teaming toolkit. Chapter 6 evaluates the toolkit on selected LLMs. Finally, Chapter 7 summarises the results and outlines directions for future work.

Chapter 2

Background on Large Language Models and Safety Risks

This chapter provides the technical and empirical background necessary to understand safety risks and adversarial behaviour in large language models (LLMs).

We first cover the architecture and training paradigm of modern LLMs and the alignment methods applied to improve their safety.

Subsequently, we survey the primary categories of safety risks and known vulnerability classes arising from LLM deployment.

2.1 Large Language Models

Large language models (LLMs) are sequence-to-sequence (encoder-decoder) or autoregressive (decoder) models based on the Transformer architecture [28].

They are trained on massive, internet-scale corpora using a self-supervised objective such as next-token prediction. After pretraining, many high-performance models undergo instruction fine-tuning and additional alignment to become more useful and safer in downstream use [31].

By 2025, publicly available and research-grade LLMs often exceed tens to hundreds of billions of parameters [12], and many incorporate advanced architectural and inference-optimisation techniques such as mixture-of-experts layers [6], retrieval-augmented generation [2], or quantisation-aware training for efficient deployment [4].

Because of the scale of their training data and capacity, LLMs encode a wide variety of linguistic patterns, factual knowledge, biases, and behavioural priors.

This scale gives them impressive generative and reasoning capabilities, but also makes them susceptible to emergent, unintended behaviours that are not trivially predictable — including safety failures, policy evasion, and social-engineering style manipulation [29].

2.1.1 Alignment and Safety Interventions

To mitigate risks from raw pre-trained models, developers commonly apply a pipeline of alignment techniques.

First, supervised fine-tuning (SFT) on curated instruction-following datasets helps the model adhere to desired task formats [31, 27].

Subsequently, reinforcement learning from human feedback (RLHF) is often used: human annotators rate model outputs, a reward model is trained on those ratings, and a

policy optimisation algorithm such as Proximal Policy Optimisation (PPO) updates the model to align it to human preferences [17, 25].

Variants such as Direct Preference Optimisation (DPO) or RLAI (reinforcement learning from AI feedback) aim to improve efficiency and scalability of alignment without sacrificing safety or quality [23, 9].

Despite alignment efforts, evidence demonstrates that safety training is not foolproof. Behavioural failures persist, especially under adversarial or adversary-chosen inputs: prompt-based attacks and jailbreaks remain a major vulnerability class (cf. [29, 33, 21]).

The fact that surface-level filtering and instruction tuning can be bypassed indicates that LLMs continue to rely on shallow heuristics rather than robust semantic safety guarantees [29, 33, 21].

2.2 LLM Safety Risks

LLM misuse and unintended outputs pose a broad array of risks.

The following taxonomy, commonly adopted in recent red-teaming efforts [22, 13], captures the primary threat vectors:

- **Malicious use:** generation of instructions or actionable content facilitating wrongdoing (e.g., construction of weapons or explosives, synthesis of illicit substances, or cyber intrusion).
- **Harassment, hate, and discrimination:** outputs that demean, harass or promote bias against individuals or protected groups.
- **Misinformation:** harmful or misleading claims that may influence user behaviour or beliefs.
- **Hallucinations:** fabricated or unfounded statements presented with confidence.
- **Privacy violations:** unintended leakage of sensitive information, either from memorised training data or through malicious prompts.
- **Self-harm / dangerous content:** generation of content promoting self-harm, suicide, or exploitation (especially of minors).
- **Emergent misuse and behavioural failures:** including instruction-following failures, refusal evasion (the model ignoring safety instructions), social-engineering exploitation, or covert manipulation over multiple turns.
- **Tool misuse / unsafe actions:** harmful commands or unintended actions when LLMs interface with external tools (e.g., code execution, browser control, or file manipulation).

These categories correspond to those used in large-scale safety benchmarks and red-team evaluation suites.

For example, frameworks such as *HarmBench* operationalise a broad set of harm categories and provide standard test sets and evaluation protocols across LLMs and red-teaming methods [13].

2.2.1 Application- and System-level Risks: OWASP LLM Top 10 (2025)

Beyond model-level behavioural vulnerabilities, real-world LLM deployments introduce additional system and integration risks. The OWASP LLM Top 10 (2025) [19] provides an industry-oriented overview of common failure modes encountered in practical LLM-based systems. For completeness, the full list is provided below, together with indicative categories reflecting the level at which each risk typically manifests:

- **LLM01: Prompt Injection** — crafted inputs override or manipulate system instructions. (*Model-level*)
- **LLM02: Sensitive Information Disclosure** — unintended leakage of private or memorised data. (*Model-level*)
- **LLM03: Supply Chain Vulnerabilities** — risks arising from compromised datasets, model weights, or third-party components. (*Deployment-level*)
- **LLM04: Data and Model Poisoning** — malicious training or fine-tuning data influencing model behaviour. (*Training-level*)
- **LLM05: Improper Output Handling** — insufficient sanitisation or validation of model outputs. (*Application-level*)
- **LLM06: Excessive Agency** — autonomous or uncontrolled agent actions producing unintended effects. (*Agent-level*)
- **LLM07: System Prompt Leakage** — extraction of hidden or system-level instructions. (*Model-level*)
- **LLM08: Vector and Embedding Weaknesses** — vulnerabilities in retrieval-augmented generation (RAG) pipelines or embedding stores. (*RAG-level*)
- **LLM09: Misinformation** — confident but incorrect or misleading outputs. (*Model-level*)
- **LLM10: Unbounded Consumption** — excessive resource usage leading to cost escalation or denial-of-service. (*Operational-level*)

Chapter 3

Red-Teaming of Large Language Models

This chapter focuses on red-teaming as a systematic methodology for evaluating the safety and robustness of large language models.

We first introduce red-teaming in the context of LLMs, including its definitions, goals, and a taxonomy of testing approaches used in red-teaming.

Subsequently, we describe common attack techniques (single-turn, multi-turn, prompt injection, universal triggers, trojans/backdoors).

Finally, we review evaluation metrics and benchmark frameworks commonly used to measure safety failures, and survey existing open-source tools and frameworks for LLM red-teaming.

The chapter concludes with a synthesis of gaps in current tooling and evaluation practice, motivating the design of the toolkit proposed in this thesis.

3.1 Red-teaming: Definitions and Methodology

Originally developed in cybersecurity and military contexts, *red-teaming* refers to the practice of simulating adversaries to probe system vulnerabilities before deployment [22].

In the context of LLMs, red-teaming denotes systematic probing of model behaviour via intentionally crafted adversarial prompts or inputs designed to circumvent safety and alignment measures, with the goal of discovering previously unknown failure modes, measuring their prevalence, and informing robust defences [24, 1].

Red-teaming campaigns in LLM settings typically have several objectives [22, 24]:

- **Discovery:** reveal novel, previously undocumented failure modes (e.g., new jailbreak styles, multi-turn attack vectors, prompt injection in application contexts).
- **Quantification:** estimate how often a model fails under adversarial conditions, enabling comparison across models and defence strategies.
- **Reproducibility:** produce repeatable test cases and evaluation pipelines so that safety regressions can be detected over time.
- **Defence hardening:** feed findings back into model tuning, safety filters, or deployment guardrails to reduce vulnerability.

Unlike static benchmarks, red-teaming targets the worst-case behaviour of a model rather than its average-case performance, focusing on adversarial inputs specifically crafted to expose vulnerabilities [13, 24]. Modern red-teaming approaches in LLMs span from fully manual adversarial testing to automated pipelines that integrate adversarial prompt generation, execution, and evaluation.

The following subsections describe manual, rule-based, and automated red-teaming methodologies in more detail.

3.1.1 Manual, Rule-based and Automated Red-teaming

Manual red-teaming: human experts write adversarial prompts, simulate realistic misuse scenarios, and attempt to elicit harmful or policy-violating responses [22]. This approach benefits from human creativity and insight into real-world misuse, including social-engineering, context-aware manipulation, and subtle or ambiguous scenarios. However, it is expensive, time-consuming, and often non-reproducible.

Rule-based / template-based testing: uses curated prompt templates (e.g., standard jailbreaks, role-play prompts, obfuscation methods, encoded instructions) or transformation rules to generate adversarial inputs systematically [24, 13]. This method is reproducible and simple, but often limited to discovering known failure classes, and does not generalise to novel or adaptive attacks.

Automated red-teaming: algorithmic or model-in-the-loop generation of adversarial prompts, for example via search, optimisation, or by using another LLM as the attacker [24, 1]. Notable work in this space includes *Tree of Attacks (TAP)* which automatically generates jailbreak prompts against black-box LLMs by iteratively refining candidates in a tree-structured search and pruning unlikely attack paths [14]. Automated methods scale well, explore large prompt spaces, and can discover novel failures, but may also generate unrealistic or trivial prompts, and often depend on the quality of the attacker or judge model and the underlying search strategy.

Automated red-teaming has recently been re-envisioned as a sequential, multi-turn process rather than isolated single-turn attempts. For example, recent work models red-teaming as a Markov Decision Process (MDP), using hierarchical reinforcement learning (RL) to optimise long-horizon attacks over entire dialogue trajectories [1].

3.2 Attack Taxonomy

Adversarial attacks against LLMs span a diverse space of techniques that exploit different aspects of model behaviour and prompting dynamics. Organising these attacks into a taxonomy is useful for understanding their mechanisms, comparing red-teaming methods, and designing systematic evaluation pipelines [22].

This section surveys the primary attack classes considered in contemporary LLM red-teaming: single-turn jailbreaks, prompt-injection attacks, multi-turn conversational manipulation, universal or transferable triggers, and trojan or backdoor-style vulnerabilities [24].

3.2.1 Single-turn Jailbreaking

Single-turn jailbreaks are adversarial prompts supplied in one-shot (single user message) that instruct the model to ignore its safety filters or system instructions [29].

Common techniques include role-play (“Pretend you are ...”), direct overrides (“Ignore previous instructions ...”), encoding or cipher-based obfuscation to hide disallowed instructions, and other prompt-engineering strategies [21].

Despite alignment training, many models remain vulnerable to these attacks [29, 33].

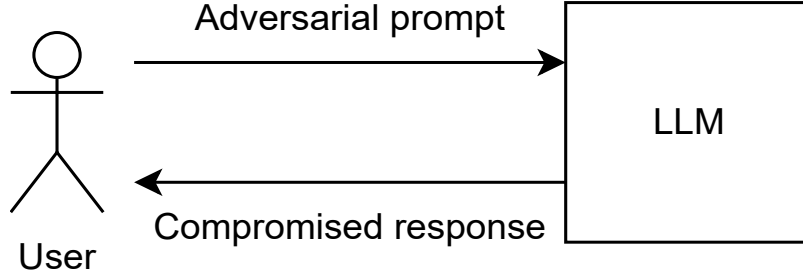


Figure 3.1: Example of a single-turn jailbreak prompt that elicits disallowed content from an LLM in one prompt. Adapted from Zhao and Zhang [32].

3.2.2 Prompt Injection

Prompt injection refers to attacks where user-supplied input (or external content, in the case of integrated applications) is directly interpreted by the LLM as instructions, potentially overriding or modifying hidden system prompts [21]. This is a major concern for real-world applications embedding LLMs (chatbots, agents, document processors, pipelines), because attackers can inject maliciously crafted instructions that the model interprets as legitimate. Empirical studies have demonstrated that many deployed LLM-based applications are vulnerable to prompt injection — for example via the *Hou Yi* attack, which compromised dozens of applications in a black-box setting [11]. Other work has shown that even automated, universal prompt injection attacks remain effective under defensive measures [10].

Prompt injection remains among the most significant security threats for LLM-based systems, as acknowledged by security guidance frameworks and cheat sheets (e.g., from OWASP) tailored to LLM applications [18, 19].

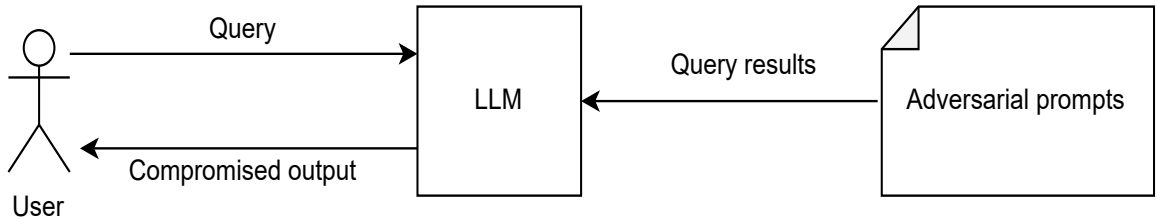


Figure 3.2: Example of an indirect prompt injection attack: a user issues a benign query, the LLM retrieves external content containing embedded instructions, and these instructions override system or developer prompts, leading to compromised output. Adapted from Greshake et al. [7].

3.2.3 Multi-turn and Conversational Attacks

Rather than achieve a jailbreak in a single prompt, adversaries may perform a sequence of manipulative steps — gradually steering the conversation, exploiting context persistence, memory, and the model’s inability to consistently refuse undesirable requests [29].

Multi-turn attacks may involve context poisoning, social-engineering style dialogue, bait-and-switch tactics, or incremental obfuscation. This vector is increasingly recognised as one of the most dangerous and under-evaluated, as multi-step interactions more closely resemble realistic misuse scenarios.

Recent research recasts automated red-teaming as a multi-turn optimisation problem — better capturing long-horizon adversarial strategies — and shows that RL-based red-teaming significantly outperforms single-turn methods in eliciting harmful content [1].

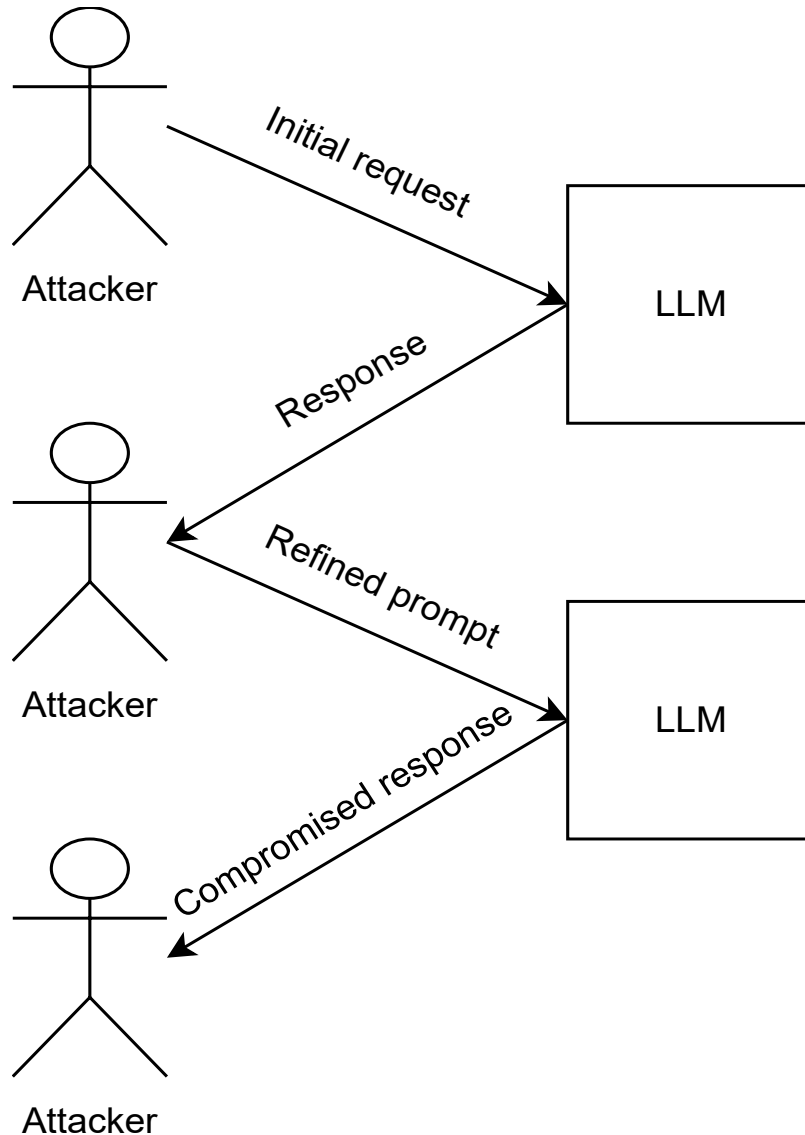


Figure 3.3: Illustration of a multi-turn jailbreak attack, in which an attacker iteratively refines prompts based on previous model responses to progressively elicit disallowed behavior. Adapted from Zhao and Zhang [32].

3.2.4 Universal / Transferable Triggers and Trojan-style Attacks

Universal triggers or transferable adversarial prompts aim to find short token sequences (prefixes, suffixes, or embedded instructions) that reliably trigger undesired model behaviour across different inputs and even across different models [33]. This makes the attack highly reusable and dangerous.

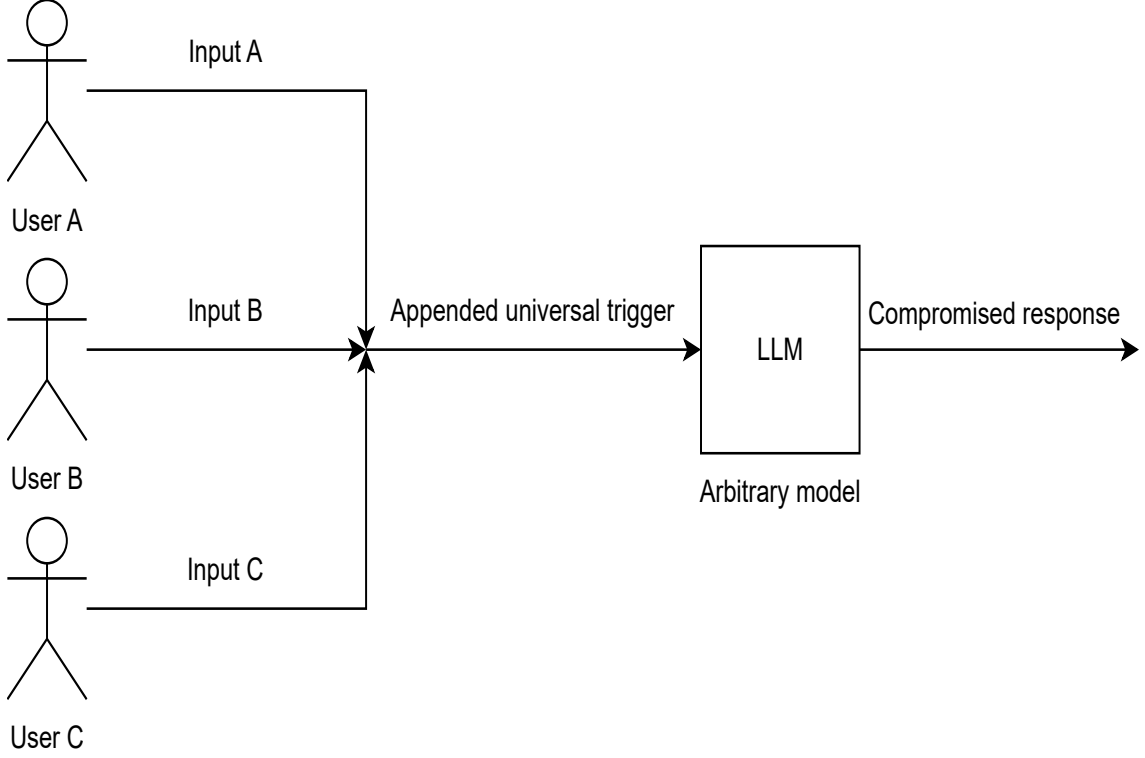


Figure 3.4: Illustration of a universal trigger attack, where a fixed trigger sequence appended to user inputs induces undesired behavior in an unmodified language model across different prompts and, in some cases, across different models.

A related but distinct class of attacks are Trojan-style (backdoor) attacks, where the model itself is modified during training or fine-tuning to respond maliciously to a specific trigger. An example is the black-box Trojan prompt attack framework *TrojLLM*, which demonstrates that stealthy trigger patterns can be embedded into widely used LLM architectures and APIs, enabling malicious behavior when the trigger is present [30]. Because the trigger is designed to appear benign, such attacks may bypass input sanitization and heuristic defenses.

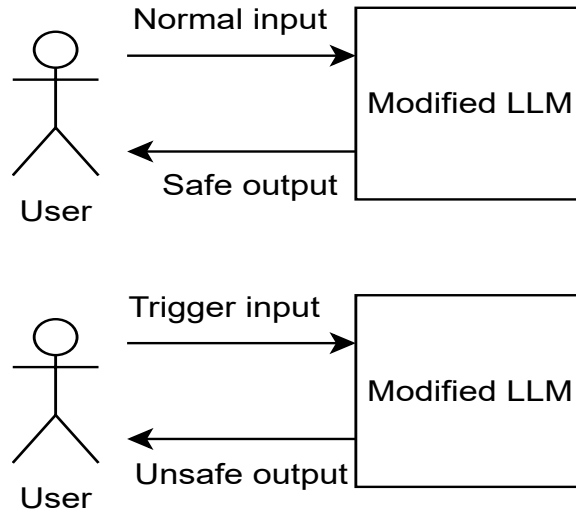


Figure 3.5: Illustration of a Trojan-style (backdoor) attack on a large language model. The modified model behaves normally for benign inputs but produces malicious or policy-violating outputs when a specific trigger pattern is present. Unlike universal trigger attacks, the vulnerability is embedded in the model during training or fine-tuning rather than in the input prompt.

3.2.5 Prompt-based Evasion and Adversarial Perturbations

Beyond explicit instructions or injected content, adversaries may attempt adversarial prompting via subtle token-level perturbations, insertions, deletions, or encodings — minimal but adversarial changes that remain semantically similar to benign prompts yet cause the model to deviate [33].

Defences based on sanitisation or heuristic filtering often fail to detect such manipulations, especially in the presence of context sensitivity or long prompts. Certified-safety approaches have been proposed to mitigate this class: for example, erase-and-check, which systematically removes tokens and reruns safety filters to detect adversarial prompt manipulations, providing (under assumptions) a safety guarantee against insertions, suffixes, or adversarial infusions up to a bounded size [8]. While promising, these methods are computationally expensive and may degrade user experience or model utility.

Because these perturbation-based attacks operate at the token level and can be obfuscated, they represent a difficult-to-detect threat, especially when combined with other attack modalities (multi-turn, universal triggers, trojans).

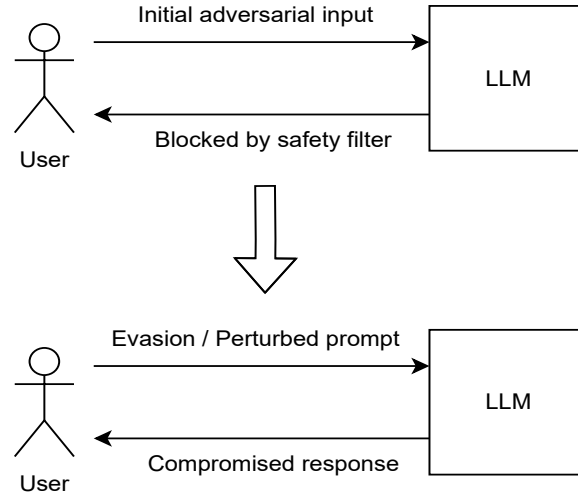


Figure 3.6: Illustration of a prompt-based evasion attack, where an attacker modifies a disallowed request using linguistic obfuscation or adversarial perturbations to bypass safety mechanisms.

3.3 Evaluation Metrics and Benchmarks

Robust and systematic evaluation is critical for red-teaming, as it enables meaningful comparison of models, attack strategies, and defences. The following metrics and benchmarks are widely adopted in contemporary LLM safety research and practice [22].

3.3.1 Quantitative Metrics

- **Attack Success Rate (ASR):** the proportion of adversarial attempts (prompts) that succeed in eliciting harmful or policy-violating outputs. ASR is the most common metric, but it depends heavily on the definition of „harmful“ and on the quality of the judge (classifier, LLM-judge, human) [29, 13, 24].
- **Refusal Bypass Rate:** a variant of ASR that measures how often safe-mode refusals or safety filters are bypassed — i.e., the model produces a disallowed output rather than refusing the request [29, 22].
- **Judge Reliability Metrics:** when using automated judges (regex, classifier, LLM-based), it is crucial to measure false positives / false negatives, calibration error, and inter-annotator agreement (if doing human calibration). Poor judges can substantially distort ASR and other metrics [1, 22].
- **Robustness Metrics:** measure the stability of attacks or defences under input variations — such as paraphrasing, prompt reordering, model temperature differences, different seeds (randomness), context shuffling, or small perturbations [13, 14].
- **Transferability / Generality:** measure how well attacks discovered on one model or configuration transfer to other models, prompts, or deployments [33, 13].
- **Conversational / Trajectory Metrics:** in multi-turn attack scenarios, metrics may capture success over a dialogue trajectory (e.g., whether harmful content emerges at any point), time-to-failure, or complexity (number of turns required) [24, 14].

3.3.2 Benchmarks and Standardised Suites

To enable systematic comparison across red-teaming methods and LLMs, several standard benchmarks have recently been developed. A prominent example is *HarmBench*, which provides an open-source evaluation framework, a large pool of red-teaming methods and a diverse set of target models and defences [13].

HarmBench enables reproducible large-scale red-team evaluation and supports both attack- and defence-side experiments. Its release marks a milestone toward standardising LLM safety evaluation pipelines.

Despite this progress, many prior works still rely on ad-hoc prompt sets, non-public prompt libraries, or private internal red-team pipelines. In addition, while single-turn attack benchmarks are relatively common, multi-turn or conversational red-teaming remains underrepresented in publicly available benchmark suites.

3.4 Open-source Red-Teaming Tools and Frameworks

A variety of open-source frameworks and research prototypes have been developed to support systematic red-teaming of LLMs. Their design choices, capabilities, and limitations differ significantly, reflecting diverse goals and use cases.

3.4.1 Frameworks for rule-based or template-based red-teaming

Frameworks such as *PyRIT* provide a modular architecture for plugging in model backends, defining prompt templates, logging results, and running interactive or batch red-teaming sessions. These tools are often lightweight and well-suited for smaller projects or smaller compute budgets, but tend to lack advanced automated search or generation capabilities [15].

Another widely used toolkit is *garak*. *garak* provides „probes, generators, and detectors“: probes manage attack logic; generators abstract target models (LLMs, dialog systems, or any component taking text and returning text); detectors assess whether output indicates a successful attack; and the framework compiles results into human-readable reports (HTML and JSON). This design allows red-teaming across a variety of model backends and output modalities, but — as with other rule-based frameworks — it may not support automated, learning-based attack generation or multi-turn conversational scenarios out of the box [3].

3.4.2 Automated Red-Teaming Frameworks

More advanced frameworks attempt to integrate attack generation, execution, and evaluation into a unified pipeline. Notable recent work includes *MAD-MAX*, a modular adversarial red-teaming framework designed to allow multiple attack strategies (template-based, search-based, LLM-driven) within a pluggable architecture [24]. Its modular nature makes it flexible and extensible, but in practice integrating it with diverse model runtimes (local LLaMA variants, API-based models, multi-modal models) and scaling up to large-scale red-teaming campaigns remains challenging due to compute demands, backend model compatibility, and evaluation infrastructure requirements.

Beyond modular frameworks, recent work explores fully automated, trajectory-based red-teaming that recasts red-teaming as a sequential decision-making process over entire dialogues. For example, the automated red-teaming approach by Belaire et al. (2025) formalises multi-turn red-teaming as a Markov Decision Process (MDP), enabling attack policies to optimise over entire conversation trajectories rather than single messages [1]. This methodology captures realistic adversarial behaviour and reveals vulnerabilities that single-turn or template-based attacks may miss.

3.4.3 Defence-oriented and Certified Safety Approaches

In response to the growing sophistication of attacks, some work focuses on hardening LLMs against adversarial prompting. For example, the framework „erase-and-check“ provides a method for certifying safety against adversarial prompt modifications (suffix insertion, insertion at arbitrary positions, adversarial infusions) under a bounded adversarial size. This method recomputes safety classification after systematically removing tokens from the prompt and offers provable safety guarantees under certain assumptions [8]. While promising, such approaches are often computationally expensive and may impair user experience or model usability.

Moreover, the existence of Trojan prompt attacks (as demonstrated by frameworks such as *TrojLLM*) indicates that even seemingly benign prompts may embed stealthy triggers that cause harmful behaviour, complicating defence strategies and necessitating robust input sanitisation, runtime monitoring, or adversarial-resistant prompt encoding [30].

To summarise the capabilities, design choices, and limitations of the reviewed frameworks, tables 3.1 and 3.2 provide a comparison of rule-based, prompt-injection, automated, and certified-safety tools.

Framework	Attack types	Backends	Multi-turn	Key notes
PyRIT [15]	Rule/template prompts, basic injection.	API-based (OpenAI, Anthropic); limited local support.	Partial	Lightweight; regex/classifier judges; limited automation.
garak [3]	Template attacks, probes, obfuscation.	API + HuggingFace; extensible generators.	Partial	Regex/classifier detectors; widely used; limited multi-turn.
TAP [14]	Search-based jailbreaks.	Black-box APIs; local wrappers.	No	Query-efficient; single-shot focus.
Prompt-injection frameworks (e.g., HouYi) [11]	Direct/indirect injection; document-based attacks.	LLM-integrated applications.	Yes	Realistic workflows; task-specific checks.

Table 3.1: Comparison of rule-based and prompt-injection frameworks.

Framework	Attack types	Backends	Multi-turn	Key notes
MAD-MAX [24]	Template/LLM/search attacks.	API + local via plugins.	Yes	Modular; multi-turn; high compute cost.
Automatic LLM Red-teaming (Belaire et al.) [1]	LLM-as-attacker; RL optimisation; trajectory attacks.	API/local (model-in-loop).	Yes	Trajectory MDP; powerful but expensive.
Erase-and-Check [8]	Token deletions, suffix detection (defence).	Model-agnostic inputs.	N/A	Provable guarantees; computationally heavy.

Table 3.2: Comparison of automated and certified safety frameworks.

The contrasts highlighted in these tables reveal several systematic limitations across current tooling, which motivate the research gaps discussed in the following section.

3.5 Synthesis and Research Gap

The literature and tools surveyed above show that, although the research community has developed a rich taxonomy of harms, attack strategies, and evaluation metrics, substantial gaps remain that limit the effectiveness, accessibility, and reliability of red-teaming for large language models.

- **Lack of lightweight, modular, and extensible tooling:** existing frameworks typically focus either on lightweight, rule-/template-based probing or on large-scale automated red-teaming, but rarely combine both. There is a need for toolkits that are accessible to small research teams or individual developers yet support modern LLM backends, extensible attack and evaluator plugins, and reproducible logging and benchmarking.
- **Judge reliability and evaluation consistency:** many red-teaming efforts rely on heuristic or automated judges (regex filters, simple classifiers) that often lack calibration, robustness, or reproducibility. As a result, reported Attack Success Rates (ASR) may misrepresent true safety risk. Certified-safety methods (e.g., erase-and-check) provide stronger guarantees, but are computationally expensive and may be impractical for routine use.
- **Underrepresentation of multi-turn and real-world attack vectors:** most prior work and benchmarks focus on single-shot prompts; multi-turn conversational attacks — which more accurately model real-world adversaries — remain under-evaluated. Recently proposed trajectory-based red-teaming frameworks help, but are not yet part of standard open-source toolkits.
- **Difficulty integrating across diverse model backends and deployment contexts:** LLMs are deployed in varied settings (local open-source models, API-based proprietary models, multi-modal agents, applications with external tools), yet existing tools often lack abstractions or adapters covering this diversity.
- **Limited transparency and reproducibility of prompt libraries and experiment artefacts:** many studies do not release their full prompt sets, random seeds, or evaluation logs, making independent replication or longitudinal safety regression testing difficult.

These gaps motivate the design goals for the toolkit developed in this thesis: namely, modularity and pluggability (attack generators, model adapters, evaluators), reproducible experiment manifests and logging, hybrid judging (balancing cost and fidelity), and explicit support for multi-turn conversational testing. By addressing these gaps, the proposed toolkit aims to lower the barrier to entry for systematic LLM red-teaming for researchers, students, and small teams.

Chapter 4

Design

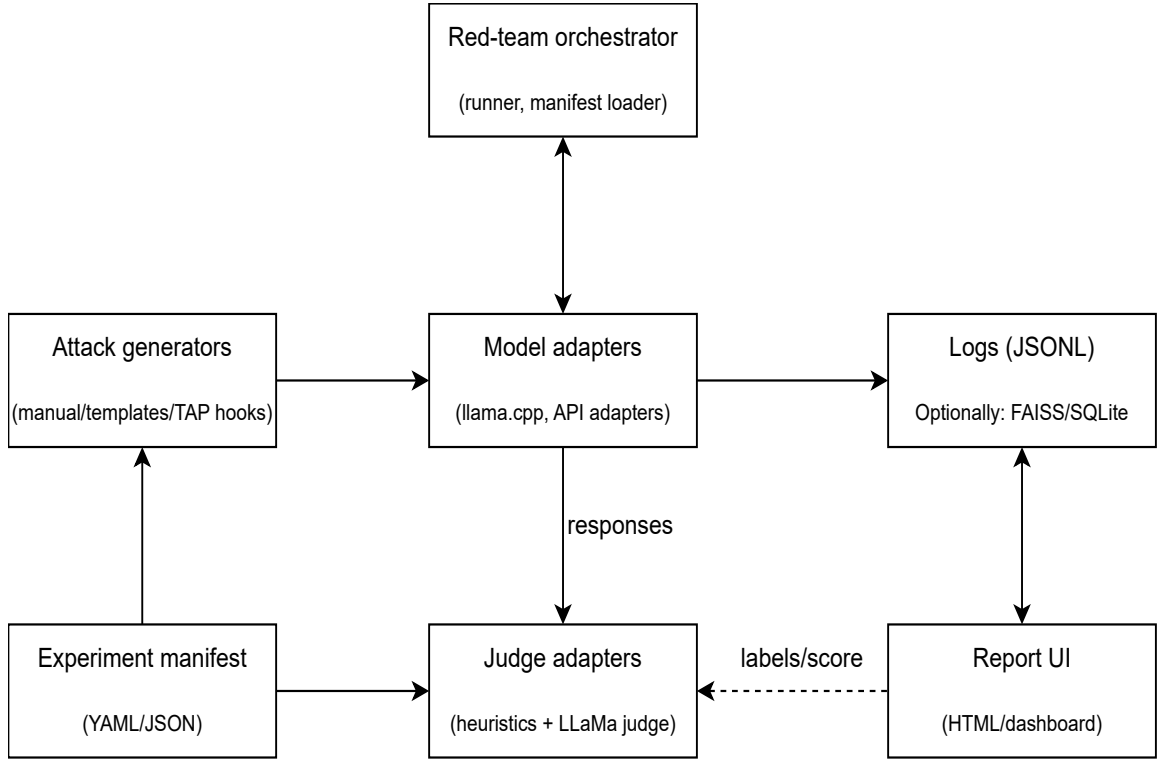


Figure 4.1: High-level architecture of the proposed red-teaming toolkit. The system is organised around a central orchestrator that coordinates attack generators, model adapters, and a hybrid judging pipeline. Experiments are defined declaratively via manifests and logged in a reproducible JSONL format, with optional indexing and reporting via a lightweight UI.

4.1 Requirements and Goals

- Functional requirements (attack generators, model adapters, evaluators, logging).
- Non-functional requirements (modularity, low compute footprint, reproducibility, extensibility).

- Mapping to gaps identified in Chapter 2.

4.2 High-level Architecture

- Component diagram and narrative.
- Data and control flows.

4.3 Module Specifications

4.3.1 Model Adapters

4.3.2 Attack Generators

4.3.3 Judges and Evaluation

4.3.4 Experiment Manifests and Logging

4.3.5 Plugin System and Extensibility

4.4 Default Test Suite and OWASP Mapping

- Default test-cases, how OWASP Top-10 items map to tests.

4.5 Security, Privacy and Ethical Considerations

- Rate-limits, safe defaults, handling PII, red-team ethics.

4.6 Proof-of-Concept Deployment

- Minimal runtime diagram (local vs API backends), tools and tech choices.

Chapter 5

Implementation

Chapter 6

Evaluation

Chapter 7

Conclusion

Bibliography

- [1] BELAIRE, R.; SINHA, A. and VARAKANTHAM, P. *Automatic LLM Red Teaming*. 2025. Available at: <https://arxiv.org/abs/2508.04451>.
- [2] BORGEAUD, S.; MENSCH, A.; HOFFMANN, J.; CAI, T.; RUTHERFORD, E. et al. *Improving language models by retrieving from trillions of tokens*. 2022. Available at: <https://arxiv.org/abs/2112.04426>.
- [3] DERCZYNSKI, L.; GALINKIN, E.; MARTIN, J.; MAJUMDAR, S. and INIE, N. *Garak: A Framework for Security Probing Large Language Models*. 2024. Available at: <https://arxiv.org/abs/2406.11036>.
- [4] DETTMERS, T.; PAGNONI, A.; HOLTZMAN, A. and ZETTLEMOYER, L. *QLoRA: Efficient Finetuning of Quantized LLMs*. 2023. Available at: <https://arxiv.org/abs/2305.14314>.
- [5] ESMAILZADEH, Y. *Potential Risks of ChatGPT: Implications for Counterterrorism and International Security*. 2023. Available at: https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4461195.
- [6] FEDUS, W.; ZOPH, B. and SHAZEER, N. *Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity*. 2022. Available at: <https://arxiv.org/abs/2101.03961>.
- [7] GRESHAKE, K.; ABDELNABI, S.; MISHRA, S.; ENDRES, C.; HOLZ, T. et al. *Not what you’ve signed up for: Compromising Real-World LLM-Integrated Applications with Indirect Prompt Injection*. 2023. Available at: <https://arxiv.org/abs/2302.12173>.
- [8] KUMAR, A.; AGARWAL, C.; SRINIVAS, S.; LI, A. J.; FEIZI, S. et al. *Certifying LLM Safety against Adversarial Prompting*. 2025. Available at: <https://arxiv.org/abs/2309.02705>.
- [9] LEE, H.; PHATALE, S.; MANSOOR, H.; MESNARD, T.; FERRET, J. et al. *RLAIF vs. RLHF: Scaling Reinforcement Learning from Human Feedback with AI Feedback*. 2024. Available at: <https://arxiv.org/abs/2309.00267>.
- [10] LIU, X.; YU, Z.; ZHANG, Y.; ZHANG, N. and XIAO, C. *Automatic and Universal Prompt Injection Attacks against Large Language Models*. 2024. Available at: <https://arxiv.org/abs/2403.04957>.
- [11] LIU, Y.; DENG, G.; LI, Y.; WANG, K.; WANG, Z. et al. *Prompt Injection attack against LLM-integrated Applications*. 2024. Available at: <https://arxiv.org/abs/2306.05499>.

- [12] LU, X.; LIU, Z.; LIUSIE, A.; RAINA, V.; MUDUPALLI, V. et al. *Blending Is All You Need: Cheaper, Better Alternative to Trillion-Parameters LLM*. 2024. Available at: <https://arxiv.org/abs/2401.02994>.
- [13] MAZEIKA, M.; PHAN, L.; YIN, X.; ZOU, A.; WANG, Z. et al. *HarmBench: A Standardized Evaluation Framework for Automated Red Teaming and Robust Refusal*. 2024. Available at: <https://arxiv.org/abs/2402.04249>.
- [14] MEHROTRA, A.; ZAMPETAKIS, M.; KASSIANIK, P.; NELSON, B.; ANDERSON, H. et al. *Tree of Attacks: Jailbreaking Black-Box LLMs Automatically*. 2024. Available at: <https://arxiv.org/abs/2312.02119>.
- [15] MUNOZ, G. D. L.; MINNICH, A. J.; LUTZ, R.; LUNDEEN, R.; DHEEKONDA, R. S. R. et al. *PyRIT: A Framework for Security Risk Identification and Red Teaming in Generative AI System*. 2024. Available at: <https://arxiv.org/abs/2410.02828>.
- [16] OPENAI. *Red Teaming Network*. 2023. Available at: <https://openai.com/blog/red-teaming-network>.
- [17] OUYANG, L.; WU, J.; JIANG, X.; ALMEIDA, D.; WAINWRIGHT, C. L. et al. *Training language models to follow instructions with human feedback*. 2022. Available at: <https://arxiv.org/abs/2203.02155>.
- [18] OWASP CHEAT SHEET SERIES. *LLM Prompt Injection Prevention Cheat Sheet*. 2025. Available at: https://cheatsheetseries.owasp.org/cheatsheets/LLM_Prompt_Injection_Prevention_Cheat_Sheet.html.
- [19] OWASP TOP 10 FOR LLMs. *Top 10 Risk & Mitigations for LLMs and Gen AI Apps*. 2025. Available at: <https://genai.owasp.org/llm-top-10>.
- [20] PARLIAMENT, E. and COUNCIL. *Regulation (EU) 2024/1689 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act)*. 2024. Available at: <https://eur-lex.europa.eu/eli/reg/2024/1689/oj>.
- [21] PEREZ, F. and RIBEIRO, I. *Ignore Previous Prompt: Attack Techniques For Language Models*. 2022. Available at: <https://arxiv.org/abs/2211.09527>.
- [22] PURPURA, A.; WADHWA, S.; ZYMET, J.; GUPTA, A.; LUO, A. et al. *Building Safe GenAI Applications: An End-to-End Overview of Red Teaming for Large Language Models*. 2025. Available at: <https://arxiv.org/abs/2503.01742>.
- [23] RAFAILOV, R.; SHARMA, A.; MITCHELL, E.; ERMON, S.; MANNING, C. D. et al. *Direct Preference Optimization: Your Language Model is Secretly a Reward Model*. 2024. Available at: <https://arxiv.org/abs/2305.18290>.
- [24] SCHOEPF, S.; HAMEED, M. Z.; RAWAT, A.; FRASER, K.; ZIZZO, G. et al. *MAD-MAX: Modular And Diverse Malicious Attack MiXtures for Automated LLM Red Teaming*. 2025. Available at: <https://arxiv.org/abs/2503.06253>.
- [25] SCHULMAN, J.; WOLSKI, F.; DHARIWAL, P.; RADFORD, A. and KLIMOV, O. *Proximal Policy Optimization Algorithms*. 2017. Available at: <https://arxiv.org/abs/1707.06347>.

- [26] SHAW, A.; YE, A.; KRISHNA, R. and ZHANG, A. X. *Agonistic Image Generation: Unsettling the Hegemony of Intention*. 2025. Available at: <https://arxiv.org/abs/2502.15242>.
- [27] STIENNON, N.; OUYANG, L.; WU, J.; ZIEGLER, D. M.; LOWE, R. et al. *Learning to summarize from human feedback*. 2022. Available at: <https://arxiv.org/abs/2009.01325>.
- [28] VASWANI, A.; SHAZEER, N.; PARMAR, N.; USZKOREIT, J.; JONES, L. et al. Attention Is All You Need. *CoRR*, 2017, abs/1706.03762. Available at: <http://arxiv.org/abs/1706.03762>.
- [29] WEI, A.; HAGHTALAB, N. and STEINHARDT, J. *Jailbroken: How Does LLM Safety Training Fail?* 2023. Available at: <https://arxiv.org/abs/2307.02483>.
- [30] XUE, J.; ZHENG, M.; HUA, T.; SHEN, Y.; LIU, Y. et al. *TrojLLM: A Black-box Trojan Prompt Attack on Large Language Models*. 2023. Available at: <https://arxiv.org/abs/2306.06815>.
- [31] ZHANG, B.; LIU, Z.; CHERRY, C. and FIRAT, O. *When Scaling Meets LLM Finetuning: The Effect of Data, Model and Finetuning Method*. 2024. Available at: <https://arxiv.org/abs/2402.17193>.
- [32] ZHAO, Y. and ZHANG, Y. *Siren: A Learning-Based Multi-Turn Attack Framework for Simulating Real-World Human Jailbreak Behaviors*. 2025. Available at: <https://arxiv.org/abs/2501.14250>.
- [33] ZOU, A.; WANG, Z.; CARLINI, N.; NASR, M.; KOLTER, J. Z. et al. *Universal and Transferable Adversarial Attacks on Aligned Language Models*. 2023. Available at: <https://llm-attacks.org>.