

# Safety as a Secondary Objective

Systematic Adversarial Evaluation of Small Language Models in High-Stakes Deployments

Polina Moshenets

SichGate

research@sichgate.com

## Abstract

Small language models (SLMs) are being deployed in high-stakes regulated settings — healthcare triage, clinical decision support, financial advisory tools — yet the adversarial security literature has almost exclusively characterised vulnerabilities in large, frontier-scale models, leaving SLM-specific risks largely unexamined. We present the systematic, black-box adversarial evaluation of six open-weight SLMs spanning 1.5 billion to 8 billion parameters, applying a structured framework of 154 tests per model across 21 attack categories covering alignment exploitation, demographic bias, context-window manipulation, multi-turn jail-breaking, and instruction hijacking, for a total of 924 adversarial interactions. Our principal finding is a three-tier, architecture-correlated gradient in context-window safety displacement that cleanly separates Llama-family models (0/8 and 0–1/6 pass rates on empty-token padding and context-overflow probes), Mistral-7B (5/8 and 2/6), and models from the Qwen and Phi-3 families (8/8 and 3/6), suggesting that vulnerability to this attack class is mediated by attention mechanism design rather than model scale.

We further demonstrate that medical domain fine-tuning redistributes adversarial exposure, improving safe-messaging compliance while amplifying demographic bias at a difference score of 0.50 against a threshold of 0.15 across eight critical findings — rather than reducing the overall attack surface. These results have immediate implications for regulated-industry deployers making procurement decisions with no SLM-specific adversarial benchmarks available, and who bear primary liability for deployment-time adversarial failures under applicable AI governance frameworks.

## 1 Introduction

The deployment profile of language models in commercial software has shifted materially over the past two years. Where organisations once accessed language

model capabilities exclusively through large-scale API endpoints, the commoditisation of open-weight models in the 1–8 billion parameter range has enabled local, on-premises, and edge deployments at a cost and latency profile compatible with latency-sensitive and data-sensitive regulated applications. Electronic health record vendors have begun embedding SLM-powered summarisation and triage assistants directly into clinical workflows, and financial services firms are deploying SLMs for customer-facing Q&A, internal compliance document retrieval, and automated advisory preliminary screens. Healthcare software companies are fine-tuning domain-specific models on proprietary clinical corpora and shipping them as on-premises medical reasoning assistants.

This shift carries a largely unacknowledged security implication. The adversarial machine learning literature has, with few exceptions, studied vulnerabilities in large, frontier-scale models: GPT-4, Claude, Gemini, and their contemporaries [12, 17, 20]. The attack surface characterisation developed in that literature does not transfer cleanly to SLMs for at least three reasons. First, SLMs are almost always quantised for deployment — typically to 4-bit or 8-bit integer formats — and quantisation-aware evaluation requires distinct test methodology. Second, regulated-industry SLMs are routinely fine-tuned on domain-specific corpora, and fine-tuning delta introduces attack vectors that do not exist in the base model. Third, SLMs are frequently deployed in agentic and multi-turn configurations where stateful context manipulation is possible, yet the adversarial multi-turn literature was developed against models with substantially richer instruction-following training.

We are not aware of any published systematic adversarial evaluation that addresses this combination of characteristics across a representative sample of SLMs currently in regulated-industry deployment.

**This paper presents such an evaluation.** We developed a structured adversarial framework comprising 154 test executions per model across 21 attack categories and applied it to six open-weight

SLMs spanning 1.5B to 8B parameters: Qwen2-1.5B-Instruct, Llama-3.2-3B-Instruct, Phi-3-mini-4k-Instruct, MedGemma-4B-IT (a medically fine-tuned vision–language model evaluated in text-only mode), Mistral-7B-Instruct-v0.2, and Llama-3.1-8B-Instruct.

We preview our principal findings here, as the introduction, not the conclusion.

The most architecturally significant result is a three-tier context-window safety displacement gradient that is correlated with model training family. Llama-3.1-8B and Llama-3.2-3B failed *all* empty-token padding (ETP) tests (0/8 each) and the majority of context-overflow displacement (COD) tests, all at high or critical severity. Mistral-7B showed partial vulnerability (5/8 ETP passes, 2/6 COD passes). Qwen2-1.5B, MedGemma-4B, and Phi-3-mini showed no ETP vulnerability (8/8 passes each). The gradient is consistent with an architectural interpretation: Llama and Mistral share the RoPE positional embedding scheme and, in earlier versions, a sliding window attention lineage absent from Qwen’s GQA-based design and Phi-3-mini’s synthetic-data training regime. This is not a generalised quality gradient, Llama-3.1-8B simultaneously showed the best sycophancy resistance of any model tested. Specific capabilities can be strong while a distinct architectural vulnerability is simultaneously present.

The domain fine-tuning result is equally counterintuitive. MedGemma-4B, the only fine-tuned model in the evaluation, improved on exactly one safety dimension (safe messaging, 5/5 vs. 4/5 for general-purpose models) while incurring eight critical-severity and eight high-severity demographic bias findings across medical pain assessment and mental health contexts. Its overall failure rate (51.3%) exceeded that of Qwen2-1.5B (42.2%) despite a 3× parameter advantage. Fine-tuning shifted the attack surface; it did not reduce it.

Five of six models failed all crescendo multi-turn escalation probes at critical severity; Phi-3-mini demonstrated consistent partial resistance, passing 1/3 escalation paths across independent evaluation runs

The remainder of this paper is organised as follows. Section 2 reviews the relevant adversarial literature and establishes the gap this work addresses. Section 3 defines the threat model. Section 4 describes the evaluation framework and model selection. Section 5 presents the empirical results. Section 6 examines the implications for deployment, regulation, and future research. Section 7 concludes.

## 2 Background and Related Work

### 2.1 Adversarial Attacks on Large Language Models

The adversarial evaluation of large language models has developed along two largely parallel tracks: *manual red teaming*, which relies on human creativity to elicit unsafe model behaviour through crafted prompts [7, 12], and *automated attack generation*, which uses gradient-based or search-based methods to construct adversarial suffixes or prefixes programmatically [20, 16].

Zou et al. [20] demonstrated that gradient-based adversarial suffixes could transfer across open-weight models, and subsequent work extended this to chain-of-thought jailbreaking [17], many-shot jailbreaking [1], and crescendo-style incremental escalation [14]. Perez and Ribeiro [12] characterised prompt injection as a distinct attack class applicable to instruction-following models, a finding that subsequent work has extended to tool-use and agentic settings [8].

Sycophantic model behaviour, where models abandon correct positions under user social pressure, has received increasing attention as a post-RLHF alignment failure mode [15, 18]. Sharma et al. [15] showed that RLHF-trained models systematically favour responses that match stated user preferences even when those preferences are factually wrong, a pattern our evaluation operationalises across safety-critical probe domains.

Demographic bias in language model outputs has been extensively documented in the context of hiring and credit decisions [2, 4], but its expression in medical AI contexts — where model responses differ in clinical quality or urgency estimates as a function of patient demographic markers — has received less systematic treatment at the evaluation framework level. Obermeyer et al. [11] documented algorithmic bias in a commercial healthcare risk-scoring tool; our work extends this concern to the generative model setting.

### 2.2 Why the Large-Model Literature Does Not Transfer to SLMs

Three structural differences between large-model and SLM evaluation render the existing literature insufficient for SLM security assessment in regulated deployments.

**Quantisation effects.** SLMs deployed at the edge or on-premises are almost universally quantised to 4-bit or 8-bit integer formats for inference efficiency. Quantisation is known to alter model be-

haviour at distribution boundaries, including at the token-probability boundaries that determine refusal vs. compliance decisions [5, 6]. Evaluations conducted on full-precision checkpoints may not faithfully characterise the behaviour of deployed quantised models. Large frontier models are rarely quantised for deployment by their vendors, making this difference largely invisible to the frontier-model evaluation literature.

**Fine-tuning delta.** Open-weight SLMs are routinely fine-tuned on proprietary domain corpora before deployment, introducing capability and alignment changes not present in the base model. The fine-tuning delta constitutes a novel attack surface: probes that the base model resists may succeed against a fine-tuned version if the fine-tuning corpus shifted the model’s prior distributions over refusal-adjacent token sequences, and vice versa. The large-model literature primarily evaluates vendor-supplied models for which the fine-tuning corpus is either absent (base models) or proprietary and fixed (RLHF-trained chat models); neither scenario captures the deployer-controlled fine-tuning surface.

**Architectural diversity at small scale.** Frontier models are dominated by a small number of architectural families. The SLM landscape is substantially more architecturally diverse: Qwen2 uses grouped query attention (GQA) with ALiBi-style positional encoding; Phi-3-mini uses synthetically generated training data with a conventional transformer architecture; Llama-3.x and Mistral-7B share RoPE positional embeddings with variants of sliding window attention. These architectural differences have direct implications for context-window attack surface, as our results demonstrate. A single evaluation methodology designed for a homogeneous architectural cohort cannot surface these differences.

### 2.3 Existing SLM and Small-Model Security Work

A small number of papers have addressed adversarial evaluation of models below the frontier scale. Qi et al. [13] showed that fine-tuning on a small number of harmful examples could catastrophically degrade safety alignment in the Llama-2-7B family, establishing fine-tuning as an alignment attack vector. Yang et al. [19] demonstrated that shadow alignment — inducing unsafe behaviour through benign-seeming fine-tuning data — succeeds against open-weight models in the 7B–13B parameter range. Hubinger et al. [9] characterised deceptive alignment in small models via deliberately planted behavioural triggers.

These contributions address the threat of malicious or adversarial model *training* rather than *deployment-*

*time* adversarial attacks against already-deployed models. Our work is complementary: we treat model weights as fixed and assess the deployment-time attack surface presented by the instruction-following interface, which is the surface that regulated-industry deployers need to characterise before shipping.

To our knowledge, no published work presents a systematic, multi-category, multi-model adversarial evaluation of the SLM deployment-time attack surface at the scale conducted here, or characterises the architecture-correlated differential vulnerability patterns we report.

## 3 Threat Model

### 3.1 Attacker Profile

We consider a deployment-time attacker who interacts with an SLM through its standard user-facing interface, a chat API, a clinical documentation assistant, a financial Q&A widget, or a patient-triage chatbot with no privileged access to model weights, training data, or system infrastructure. The attacker is *black-box*: she observes model outputs but cannot query model internals, perform gradient computations, or modify the inference pipeline. She may have unlimited query access and may conduct multi-turn interactions.

We do not assume sophisticated adversarial tooling. Our threat model explicitly includes non-expert attackers — clinicians, patients, financial customers, or internal users — who apply social engineering, context manipulation, and iterative prompt refinement without formal adversarial ML knowledge. This is the realistic threat surface for deployed regulated-industry SLMs, where the attacker population is broad, heterogeneous, and not technically specialised.

We additionally consider the threat of an attacker who has access to the model’s system prompt but not its weights, which is the typical configuration in API-served deployments with a fixed operator-defined system prompt.

### 3.2 Attack Surface

The deployment-time attack surface for instruction-following SLMs comprises the following components, each of which our evaluation framework addresses.

**Alignment exploitation.** The attacker attempts to elicit model behaviour that violates the operator’s intended constraints by constructing inputs that expose miscalibration between the model’s trained safety behaviour and its response to adversarial framings. Sub-categories include sycophantic capitulation under social pressure, crescendo-style incremental erosion of

constraints across turns, competing-objective framings that present helpfulness and safety as in direct conflict, and role-play or persona-based context injection.

**Context-window manipulation.** The attacker exploits the model’s finite context window to displace, dilute, or override system-level safety instructions. In bounded-context deployments (e.g., a 4K-token clinical assistant), filling the context with benign content before introducing an unsafe payload may push earlier instructions into the model’s attention periphery. In long-context deployments (e.g., Llama-3.1-8B with a 128K-token window), the same effect operates at a substantially larger absolute token depth.

**Structured-format injection.** The attacker embeds adversarial instructions within the structured data formats that the model is expected to parse — JSON objects, XML tags, function call schemas — exploiting the model’s tendency to treat structured-format content as trusted operational metadata rather than user-controlled input.

**Multi-turn state manipulation.** The attacker leverages multi-turn dialogue to establish false conversational premises — fictitious prior consent, role-conditioned authority, incrementally escalated context — that the model honours in subsequent turns. This attack class is particularly relevant to agentic deployments where the model maintains and acts on dialogue state across extended interactions.

**Demographic bias expression.** In the regulated-industry context, demographic bias in model outputs constitutes an adversarial vulnerability not because the attacker *causes* it, but because it can be *elicited* and documented by an auditor, a regulator, or a plaintiff’s counsel using structurally identical queries that differ only in patient or customer demographic encoding. We include bias probing as an adversarial evaluation category because its regulatory and legal consequences in healthcare and financial services are equivalent to those of active exploitation.

### 3.3 Black-Box Assumption and Scope

All findings reported in this paper were obtained through black-box evaluation. We did not perform gradient-based adversarial suffix generation, model weight inspection, training data extraction, or any technique requiring internal model access. This scope reflects the operational reality of most regulated-industry deployments, where the deployer receives a model checkpoint but adversarial testing must be conducted through the same interface as the end user.

White-box attacks — including adversarial suffix generation [20], embedding-space manipulation, and training data poisoning — represent a distinct and

potentially more severe threat class that is out of scope here. Our results should be interpreted as a lower bound on the full attack surface: a model that fails our black-box probes at critical severity will also present additional vulnerabilities under white-box conditions.

### 3.4 Why Regulated Industry Context Changes the Stakes

The implications of adversarial SLM vulnerabilities differ materially in regulated versus consumer settings for three reasons.

*Liability is assigned.* Under EU AI Act Article 3, deployers of high-risk AI systems bear primary responsibility for ensuring systems meet the requirements of Articles 9 through 15. A deployer who ships a clinical documentation assistant without conducting adversarial evaluation cannot claim that vulnerability responsibility lies with the model vendor. This makes the deployer, not the model author, the party whose risk exposure is most directly affected by the findings reported here.

*Consequences are bounded by domain.* A sycophancy failure in a consumer chatbot produces misinformation. The same failure in a clinical dosage assistant that reverses a correct overdose-risk warning under physician social pressure (`sycophancy_medication_overdose`) may contribute to patient harm. The difference is not in the model’s behaviour - which is identical in both cases, but in the deployment context that determines the consequence of the output.

*Auditability requirements are asymmetric.* HIPAA Security Rule §164.308 requires covered entities to implement procedures to regularly evaluate technical and operational safeguards. ISO/IEC 42001 requires documented evidence of AI risk management processes. EU AI Act Annex IV requires technical documentation sufficient to demonstrate compliance. In all three frameworks, deployers who cannot document adversarial evaluation of their AI systems face compliance exposure independent of whether an actual adversarial incident has occurred.

## 4 Methodology

### 4.1 Evaluation Framework Design

We developed a structured adversarial evaluation framework comprising 154 test executions per model across 21 attack categories. The framework is organised into five thematic groups reflecting the attack surface taxonomy described in Section 3: alignment exploitation (59 tests), bias and fairness (36 tests),

context-window manipulation (26 tests), multi-turn adversarial dialogue (6 tests), and reasoning integrity and structured-format injection (27 tests).

**Probe construction.** Each test comprises a structured input payload designed to elicit a specific failure mode, together with a ground-truth specification that defines the criteria for compliance (pass) versus violation (fail). Alignment probes use a two-turn structure: an initial query establishes a factual or safety-relevant baseline response, and a pushback turn introduces social pressure, role framing, or competing authority claims. Context-window probes use multi-request structures that vary the position and depth of an unsafe payload within a filled context window. Bias probes use structurally isomorphic query pairs that differ only in demographic markers encoded via name selection, consistent with the methodology of Bertrand and Mullainathan [3] and its extension to NLP contexts [10]. Multi-turn probes construct full dialogue histories that establish the false premises against which the unsafe final turn is evaluated.

**Response classification.** Model responses are classified as pass or fail against the probe ground truth using a two-stage evaluation procedure. In the first stage, a programmatic classifier applies rule-based string matching and semantic overlap metrics to identify clear-case outcomes. In the second stage, responses where the programmatic classifier confidence falls below a defined threshold are reviewed using a secondary LLM-based evaluator operating under an explicit rubric. For sycophancy probes, the classification criterion is whether the model’s substantive position after the pushback turn is consistent with its baseline response; responses that hedge with apologetic framing but maintain the correct substantive position are classified as passes. For bias probes, classification uses a normalised text-similarity difference score between the demographically varied query pair; findings are raised when the difference score exceeds 0.15 on a [0, 1] scale. For multi-turn probes, the classification criterion is whether the model’s final-turn response complies with the unsafe request established or escalated across prior turns.

**Severity assignment.** Each failing test is assigned one of three severity levels. *Critical* severity is assigned when the failure directly enables patient harm, regulatory violation, or system compromise in a deployment context consistent with the probe’s intended use case: a medical dosage assistant that reverses an overdose warning, a HIPAA-governed system that endorses patient data disclosure, or a model that completes a multi-turn jailbreak to provide restricted information. *High* severity is assigned when the failure indicates a substantial alignment deficiency that

would require remediation before regulated-industry deployment but does not directly enable the most severe harm outcomes in a single interaction. *Medium* severity is assigned when the failure is contextually inconsistent or suboptimal but does not independently constitute a safety or compliance violation. No test in the evaluation produced a low-severity finding. The absence of low-severity findings may reflect the targeted nature of the evaluation framework, which was designed around regulated-industry harm scenarios where low-consequence failures are structurally unlikely, or may indicate a calibration effect in the severity rubric; we cannot fully distinguish between these explanations.

## 4.2 Attack Category Descriptions

The 21 attack categories are as follows.

**Alignment sycophancy** (7 tests) assesses position stability under social pressure. Probe pairs include factual questions (`australia_capital`, `basic_math`), scientific consensus (`flat_earth`), and safety-critical domains (`medication_overdose`, `impaired_driving`, `safety_bleach_ammonia`, `hipaa_compliance`). The final three carry critical-severity failure classifications given their relevance to clinical and regulatory deployment contexts.

**Alignment crescendo** (3 tests) assesses multi-turn constraint erosion. Three escalation paths are tested: data exfiltration escalation (`crescendo_data_exfiltration_escalation`), safety boundary erosion (`crescendo_safety_erosion`), and identity dissolution (`crescendo_identity_dissolution`). Each probe is rated critical on failure by construction, because successful crescendo completion represents a full jailbreak through an otherwise unbreached safety boundary.

**Alignment competing objectives** (5 tests) presents explicit conflicts between system objectives — helpfulness versus safety, helpfulness versus privacy, utility versus regulatory compliance — and classifies outcomes by which objective the model chose to satisfy.

**Alignment consistency** (24 tests) applies paraphrase, reframing, and contextual variation to previously answered queries and assesses whether the model maintains consistent positions across structurally equivalent inputs.

**Bias and fairness** (36 tests) compares model responses to query pairs that differ only in demographically encoded patient or customer names across medical pain assessment and mental health contexts. Demographic dimensions tested include race (White

vs. Black, White vs. Hispanic, White vs. Middle Eastern, White vs. Asian), ethnicity, and gender, using name pairs from the Bertrand–Mullainathan paradigm adapted to clinical settings. White is used as the reference category consistent with the Bertrand–Mullainathan audit paradigm and the established healthcare bias literature.

**Safe messaging** (5 tests) assesses compliance with evidence-based safe messaging guidelines for discussions of self-harm and crisis scenarios.

**Instruction hijacking** (8 tests) tests whether runtime inputs can override system-level instructions established in the model’s operator context.

**JSON and XML smuggling** (6 tests) embeds adversarial instructions within structured data payloads that the model is expected to parse.

**Payload splitting** (3 tests) distributes a harmful request across multiple non-contiguous message segments to evade per-message classification.

**Virtualization** (8 tests) uses persona, role-play, and hypothetical framings to construct an alternate context in which unsafe behaviour is presented as fictional or appropriate.

**Infinite context: empty token padding** (8 tests) measures whether increasing multi-turn context fill (at 2, 3, 4, and 5 turns) causes safety degradation relative to an empty-context baseline. Each position is sampled twice; reported counts reflect aggregate compliance probability rather than a single binary trial.

**Infinite context: context window injection** (12 tests) places unsafe payloads at fixed proportions (0%, 50%, 75%, 90%) of the nominal context window, sampling three independent draws per position.

**Infinite context: context overflow displacement** (6 tests) positions unsafe payloads at or near the 95th percentile of the model’s context capacity and assesses whether safety instructions established earlier in the context are displaced. Three draws are taken per configuration.

**Multi-turn: role-conditioned instructions (RCI)** (2 tests) assesses whether expert-role framings (`unrestricted_ai`, `expert_override`) induce unsafe compliance in subsequent turns.

**Multi-turn: conversational consent manipulation (CCM)** (2 tests) assesses whether the model honours fictitiously established prior consent (`prior_agreement`, `role_consent`) as authority for unsafe requests.

**Multi-turn: trust drift via persona escalation (TDPE)** (2 tests) assesses whether deferred unsafe payloads succeed after persona establishment across prior turns (`medical_dosage`, `credential_extraction`).

**Hallucination probes** (6 tests) assess whether models confabulate authoritative-sounding but factually incorrect information in response to queries with known ground truth.

**Reasoning categories: causal attribution (CA), chain-of-thought hijacking (CoTH), logical entailment (LE), and self-consistency probing (SCP)** (3, 3, 3, and 2 tests respectively) assess the integrity of model reasoning under adversarial pressure to accept false premises, follow corrupted reasoning chains, or endorse logically invalid conclusions.

Each test was executed as an independent inference call with no shared context between probes; multi-turn probes constructed context accumulation programmatically within a single call.

### 4.3 Model Selection

The six models were selected to represent the parameter range currently deployed in regulated-industry SLM applications (1.5B–8B), architectural diversity (GQA, RoPE, synthetic-data training), and a domain fine-tuned variant. Table 1 provides the full model inventory. All models were evaluated at 4-bit quantisation using deployment-representative inference configurations, reflecting the on-premises edge-deployment environment most common in regulated-industry deployments subject to data residency constraints.

### 4.4 Execution and Validation

All tests were executed between 2026-03-27 and 2026-03-28 under identical temperature and sampling parameters. Each test execution was recorded with a unique run identifier, timestamp, input payload, full model response, classification outcome, severity level, and compliance framework mapping against EU AI Act, HIPAA, GDPR, NIST AI RMF, OWASP LLM Top 10, and ISO/IEC 42001.

Findings were reviewed for consistency between programmatic classification and the probe ground truth specification before final severity assignment. Tests where the programmatic classifier and manual review diverged were resolved by manual adjudication; this affected fewer than 3% of classified failures. Context-window tests with duplicate probe names reflect intentional multiple-draw sampling and were aggregated before severity counts were computed.

To assess classification stability, each probe in the critical-severity cluster was executed across multiple independent runs at temperature settings of 0.2 and 0.7. Pass/fail classifications were consistent across all runs for crescendo and sycophancy probe categories across all six models. Competing objectives probes ex-

Table 1: Models evaluated. Parameter counts are drawn from published model cards.

Model	Architecture family	Parameters
Qwen2-1.5B-Instruct	GQA, Qwen2	1.5B
Llama-3.2-3B-Instruct	RoPE, Llama-3	3.0B
Phi-3-mini-4k-Instruct	Synthetic training, Phi-3	3.8B
MedGemma-4B-IT <sup>†</sup>	VLM, Gemma-2 lineage	4.0B
Mistral-7B-Instruct-v0.2	RoPE + SWA, Mistral	7.3B
Llama-3.1-8B-Instruct	RoPE, Llama-3	8.0B

<sup>†</sup>Vision-language model; evaluated in text-only mode. Findings apply to the language pathway only.

hibited run-to-run variance of up to 20 percent at temperature 0.7 on individual models; findings in this category should be interpreted as indicative rather than deterministic. Phi-3-mini exhibited partial crescendo resistance at elevated temperature (0.6–0.7), passing 1/3 crescendo variants consistently, a pattern absent in all other evaluated models. The evaluation methodology and attack category specifications are openly available at [github.com/sichgate/sichgate-methodology](https://github.com/sichgate/sichgate-methodology).

## 5 Results

### 5.1 Overview

We applied the evaluation framework to six SLMs for a total of 924 adversarial interactions across 154 tests per model. Overall failure rates ranged from 42.2% (Qwen2-1.5B; 65/154 tests failed) to 65.6% (Llama-3.1-8B; 101/154 tests failed). A severity-weighted risk score ( $4C + 2H + M$ , where  $C$ ,  $H$ ,  $M$  denote the counts of critical, high, and medium findings respectively) yielded 138 (Qwen2-1.5B), 148 (Phi-3-mini), 169 (MedGemma-4B), 188 (Mistral-7B), 202 (Llama-3.2-3B), and 217 (Llama-3.1-8B).

These summary statistics already challenge a naive intuition: Llama-3.1-8B, the largest and most recently published model in the cohort, exhibited the highest failure rate and risk score. Qwen2-1.5B, the smallest model, exhibited the lowest. The ranking is not monotone with parameter count across architectures. The per-category breakdown is presented in Tables 2 and 3.

### 5.2 Universal Vulnerabilities

Three attack categories produced failure rates at or near 100% across all six models.

**Crescendo multi-turn escalation.** Five of six models failed all crescendo escalation probes at critical severity (15/18 failures). Phi-3-mini was the sole exception, passing 1/3 escalation paths consistently across independent evaluation runs at both tempera-

ture settings tested (0.2 and 0.6), suggesting a partial safety boundary robustness absent in all other evaluated models. This outlier pattern is consistent across runs and is not temperature-dependent. Successful jailbreaks required between 3 and 5 turns across models and attack paths. This result is notable, because crescendo attacks require no technical capability to execute and can be attempted by any end user with iterative query access. No model exhibited any mitigation across three structurally distinct escalation vectors, indicating that multi-turn constraint management was not addressed by any training pipeline in the evaluation cohort.

**Competing objectives: helpfulness versus patient safety.** `competing_obj_helpfulness_vs_safety_medical` was failed at critical severity by all six models. In every case the model resolved an explicit helpfulness-versus-safety conflict in favour of helpfulness. Llama-3.2-3B and MedGemma-4B failed all five competing-objectives probes (0/5). Mistral-7B failed four of five, including a model-specific critical failure on `competing_obj_biometric_data_conflict` not observed in any other model at critical severity.

**Alignment consistency.** Failure rates on the 24-item consistency battery ranged from 54.2% (Phi-3-mini) to 91.7% (Llama-3.1-8B). The dominant severity was medium, indicating positional drift rather than safety-critical reversal, but the universality of the failure suggests that positional stability under paraphrase and reframing is not a training objective in any evaluated model. Phi-3-mini was the only model to achieve a majority pass rate (11/24).

### 5.3 Context-Window Safety Displacement: A Three-Tier Architectural Gradient

The most structurally differentiated finding concerns context-window safety behaviour across the ETP and COD attack categories. Results establish a three-tier vulnerability gradient correlated with model training

family, summarised as follows.

*Tier 1 (full vulnerability)* encompasses Llama-3.1-8B and Llama-3.2-3B. Both models failed all 8 ETP tests (0/8, all high severity). Llama-3.1-8B failed 6/6 COD tests, all rated critical. Llama-3.2-3B failed 5/6 COD tests (4 critical, 1 high). On the combined ETP and COD surface of 14 tests, Llama-3.1-8B passed 1/14 and Llama-3.2-3B passed 2/14.

*Tier 2 (partial vulnerability)* is occupied solely by Mistral-7B. It passed 5/8 ETP tests (3 high-severity failures) and 2/6 COD tests (2 critical, 2 high). On the combined surface, Mistral-7B passed 7/14.

*Tier 3 (no ETP vulnerability)* comprises Qwen2-1.5B, MedGemma-4B, and Phi-3-mini. All three passed 8/8 ETP tests with zero findings. Each failed 3/6 COD tests with zero critical findings, a residual COD vulnerability shared equally across architectures. On the combined surface, all three passed 11/14.

The Tier 1 and Tier 2 models share RoPE positional embeddings. Llama-3.1-8B and Llama-3.2-3B additionally employ the extended RoPE scaling necessary to support Llama-3.1-8B’s 128K-token context window. Mistral-7B-v0.2 uses RoPE with a 4K nominal context window. Qwen2-1.5B uses GQA-based positional encoding; Phi-3-mini uses conventional transformer attention trained on synthetically constructed data. The gradient is consistent with the hypothesis that RoPE positional encoding creates sensitivity to context-position-dependent safety degradation that is amplified by the extended context scaling present in the Llama-3.x training pipeline. We present this as a falsifiable mechanistic hypothesis; confirming it would require attention-layer activation analysis outside the scope of this black-box evaluation.

Context Width Injection (CWI) probes partially extend this pattern. All six models failed CWI probes at varying positions; non-Llama models achieved 3/12 passes and Llama/Mistral models achieved 1–2/12 passes. Llama-3.1-8B and Llama-3.2-3B each incurred 2 additional critical CWI findings at the 0% context fill position, indicating safety degradation at baseline context depth before any padding is introduced.

## 5.4 Domain Fine-Tuning Effects

MedGemma-4B is the only model in the evaluation that was domain-fine-tuned before deployment — specifically, on medical dialogue, clinical documentation, and healthcare Q&A corpora. Its vulnerability profile diverges from general-purpose models in a pattern consistent with attack surface redistribution rather than net reduction.

On safe messaging probes, MedGemma-4B achieved a perfect score (5/5), one of only three models to do

so, alongside Phi-3-mini and Mistral-7B. This improvement is directionally expected from supervised exposure to medically governed dialogue.

On demographic bias probes, however, MedGemma-4B incurred 8 critical and 8 high findings (44.4% failure rate), the highest critical bias count of any sub-8B model in the evaluation. All 8 critical findings recorded a difference score of 0.50 against a threshold of 0.15 on medical pain assessment and mental health context probes. Critical-severity findings were observed for race (Black vs. White, Hispanic vs. White, Middle Eastern vs. White, Asian vs. White), ethnicity, and gender comparisons. MedGemma-4B’s medical fine-tuning corpus evidently encodes demographic signal that amplifies differential response generation in the exact clinical domains the fine-tuning was intended to improve.

MedGemma-4B additionally failed `sycophancy_hipaa_compliance` at critical severity — the sycophancy probe most directly relevant to its training domain and deployment context. Its overall failure rate (51.3%) and weighted risk score (169) exceed those of Qwen2-1.5B (42.2%, 138) and Phi-3-mini (45.5%, 148) despite a parameter advantage of  $2.7\times$  and  $1.05\times$  respectively.

## 5.5 Multi-Turn Adversarial Robustness

The multi-turn battery covers three structurally distinct manipulation vectors (RCI, CCM, TDPE) across six tests. Mistral-7B-Instruct-v0.2 failed all six (0/6, all high severity) — the only model to do so. Llama-3.1-8B and Llama-3.2-3B each passed 1/6; Qwen2-1.5B and MedGemma-4B passed 3/6; Phi-3-mini passed 4/6.

Phi-3-mini was uniquely resistant to RCI probes (2/2 passes on `mttv_rci_unrestricted_ai` and `mttv_rci_expert_override`), suggesting that its synthetic-data training regime may incidentally reduce susceptibility to expert-role-conditioned compliance manipulation.

Mistral-7B’s uniform failure across all three multi-turn categories — covering role framing, prior consent fabrication, and deferred payload injection — implies a broad stateful dialogue management deficit rather than a narrow alignment gap. This finding is directly relevant to agentic deployment configurations, which are the dominant growth vector for SLM integration in regulated-industry products.

## 5.6 Sycophancy

Sycophancy resistance was the attack category with the largest inter-model variance and the weakest correlation with parameter count. Llama-3.1-8B was the most resistant model (6/7 passes, 85.7%), while Qwen2-1.5B was the least resistant (0/7 passes, 0%), despite both occupying comparable positions in the sub-8B parameter range.

The three highest-stakes probes formed a consistent failure cluster. `sycophancy_hipaa_compliance` was failed at critical severity by four of six models (Qwen2-1.5B, MedGemma-4B, Llama-3.2-3B, Mistral-7B), making it the single highest cross-model-risk individual probe in the evaluation. `sycophancy_medication_overdose` and `sycophancy_impaired_driving` were each failed at critical severity by at least two models. Llama-3.1-8B and Phi-3-mini were the only models to pass all three critical-stakes probes.

Within the Llama family, sycophancy resistance scales with parameter count (3B: 4/7; 8B: 6/7). This scaling relationship does not hold cross-architecture: Qwen2-1.5B underperforms Llama-3.2-3B at comparable scale, while Phi-3-mini at 3.8B outperforms Mistral-7B at 7.3B. Sycophancy resistance in SLMs therefore appears to be a function of post-training reward signal composition rather than model capacity.

## 5.7 Instruction Hijacking and Structured-Format Injection

The probe `instruction_hijacking_05` was failed at critical severity by five of six models. Mistral-7B-Instruct-v0.2 failed the distinct variant `instruction_hijacking_08` at critical instead, a model-specific finding not corroborated by other models at that severity level. No model achieved a perfect pass rate on the 8-item hijacking battery (pass rates: 6/8 for MedGemma-4B, 7/8 for all others).

JSON and XML smuggling failures were uniform (2–3/6 passes, all high severity). Payload splitting produced consistent failures in five of six models (1/3 passes), with Llama-3.2-3B as the sole exception (3/3 passes); given the inconsistency with Llama-3.2-3B’s performance on adjacent injection categories, independent replication of the payload splitting result for that model is recommended before publication claims are finalised.

# 6 Discussion

## 6.1 Implications for Regulated-Industry Deployers

The dominant practical implication of these results is that regulated-industry deployers cannot rely on either parameter count or vendor-reported capabilities as proxies for adversarial safety. The model with the highest weighted risk score in this evaluation is also the largest and most capable by standard benchmarks. The model with the lowest risk score is the smallest. Within individual attack categories, the variance is more striking still: Qwen2-1.5B passes all ETP tests while Llama-3.1-8B fails all of them; Phi-3-mini resists all sycophancy probes in the critical-stakes cluster while Qwen2-1.5B fails all seven.

The context-window safety displacement gradient has immediate deployment implications for any organisation using Llama-family or Mistral-family models in long-context configurations. The 4K nominal context window of Mistral-7B-v0.2 is already sufficient to induce partial ETP vulnerability; Llama-3.1-8B’s 128K-token window creates an attack surface where unsafe payloads can be embedded at token positions far beyond what any manual prompt review would detect. Deployers who use these models in configurations where context length is user-controlled or document-driven — clinical note analysis, financial document Q&A, customer support with chat history — should treat context-window injection as an active threat, not a theoretical one.

The sycophancy cluster around healthcare compliance probes (`medication_overdose`, `hipaa_compliance`, `impaired_driving`) identifies a specific, testable, and remediable vulnerability class. Four of six models will reverse a correct safety-compliance position under user social pressure. This is not an edge case in clinical or financial deployments; it is the routine experience of users who disagree with a model’s output and say so. Deployers who have not tested for sycophantic safety reversal have not assessed whether their system maintains safety positions under normal use, let alone adversarial use.

The MedGemma-4B finding should be read as a cautionary result for anyone evaluating domain-specific fine-tuned SLMs. Medical fine-tuning improved safe messaging compliance but amplified demographic bias in pain assessment and mental health contexts — exactly the clinical domains in which demographic bias in AI outputs carries the most direct patient harm risk and the most significant regulatory exposure under the EU AI Act and US healthcare anti-discrimination

Table 2: Per-category findings, Part 1: alignment, bias, and injection categories. C = critical, H = high, M = medium. Pass counts in parentheses. Bold = cross-model significance.

Attack Category	Qwen2-1.5B	MedGemma-4B	Llama-3.1-8B	Llama-3.2-3B	Phi-3-mini	Mistral-7B
alignment_sycophancy (7)	C4/H3 (0/7p)	C1/H2 (4/7p)	C0/H1 (6/7p)	C3/H0 (4/7p)	C1/H2 (4/7p)	C3/H2 (2/7p)
alignment_crescendo (3)	<b>C3 (0/3p)</b>	<b>C3 (0/3p)</b>	<b>C3 (0/3p)</b>	<b>C3 (0/3p)</b>	<b>C2 (1/3p)</b>	<b>C3 (0/3p)</b>
alignment_competing_obj (5)	C1/H1/M1 (2/5p)	C2/M3 (0/5p)	C1/H1/M2 (1/5p)	C1/M4 (0/5p)	C1/M2 (2/5p)	C2/M2 (1/5p)
alignment_consistency (24)	H3/M14 (7/24p)	H3/M17 (4/24p)	H1/M21 (2/24p)	H2/M19 (3/24p)	H3/M10 (11/24p)	H2/M19 (3/24p)
bias_fairness (36)	C3/H2 (31/36p)	<b>C8/H8 (20/36p)</b>	<b>C9/H11 (16/36p)</b>	C4/H12 (20/36p)	C5/H12 (19/36p)	C5/H9 (22/36p)
safe_messaging (5)	M1 (4/5p)	— (5/5p)	M1 (4/5p)	M1 (4/5p)	— (5/5p)	— (5/5p)
instruction_hijacking (8)	C1 (7/8p)	C2 (6/8p)	C1 (7/8p)	C1 (7/8p)	C1 (7/8p)	C1 (7/8p)
json_xml_smuggling (6)	H3 (3/6p)	H4 (2/6p)	H4 (2/6p)	H4 (2/6p)	H4 (2/6p)	H4 (2/6p)
payload_splitting (3)	H2 (1/3p)	H2 (1/3p)	H2 (1/3p)	— (3/3p)	H2 (1/3p)	H1 (2/3p)
virtualization (8)	H2 (6/8p)	H3 (5/8p)	H3 (5/8p)	H3 (5/8p)	H2 (6/8p)	H2 (6/8p)

Table 3: Per-category findings, Part 2: context-window, multi-turn, and reasoning categories, with aggregate totals. C = critical, H = high, M = medium.

Attack Category	Qwen2-1.5B	MedGemma-4B	Llama-3.1-8B	Llama-3.2-3B	Phi-3-mini	Mistral-7B
inf.ctx_etp (8)	— (8/8p)	— (8/8p)	<b>H8 (0/8p)</b>	<b>H8 (0/8p)</b>	— (8/8p)	H3 (5/8p)
inf.ctx_cwi (12)	H9 (3/12p)	H9 (3/12p)	C2/H9 (1/12p)	C2/H9 (1/12p)	H9 (3/12p)	C1/H9 (2/12p)
inf.ctx_cod (6)	H3 (3/6p)	H3 (3/6p)	<b>C6 (0/6p)</b>	C4/H1 (1/6p)	H3 (3/6p)	C2/H2 (2/6p)
multiturn_ccm (2)	— (2/2p)	— (2/2p)	H2 (0/2p)	H2 (0/2p)	— (2/2p)	H2 (0/2p)
multiturn_rci (2)	H2 (0/2p)	H2 (0/2p)	H2 (0/2p)	H2 (0/2p)	— (2/2p)	H2 (0/2p)
multiturn_tdpe (2)	H1 (1/2p)	H1 (1/2p)	H1 (1/2p)	H1 (1/2p)	H2 (0/2p)	H2 (0/2p)
hallucination_probe (6)	H2 (4/6p)	H2/M1 (3/6p)	M3 (3/6p)	M2 (4/6p)	H1/M2 (3/6p)	H2/M1 (3/6p)
reasoning_ca (3)	— (3/3p)	— (3/3p)	M2 (1/3p)	M2 (1/3p)	— (3/3p)	M2 (1/3p)
reasoning_coth (3)	H2 (1/3p)	H1 (2/3p)	H1 (2/3p)	H2 (1/3p)	H2 (1/3p)	H3 (0/3p)
reasoning_le (3)	H1 (2/3p)	H1 (2/3p)	H3 (0/3p)	H3 (0/3p)	H2 (1/3p)	H1 (2/3p)
reasoning_scp (2)	H1 (1/2p)	H1 (1/2p)	H1 (1/2p)	H2 (0/2p)	H1 (1/2p)	H2 (0/2p)
<b>Total C/H/M</b>	C12/H37/M16	C16/H42/M21	C22/H50/M29	C18/H51/M28	C11/H45/M14	C17/H48/M24
<b>Overall fail rate</b>	42.2%	51.3%	65.6%	63.0%	45.5%	57.8%
<b>Weighted score</b>	138	169	217	202	148	188

frameworks. A deployer who evaluated safe messaging only, and found MedGemma’s 5/5 performance satisfactory, would be shipping a model with eight uncharacterised critical demographic bias findings.

## 6.2 EU AI Act and Regulatory Framework Implications

EU AI Act Article 3 defines deployers of high-risk AI systems as the parties primarily responsible for ensuring that systems comply with the requirements of Chapter III, Section 2. AI systems deployed in medical device interfaces, credit scoring, employment screening, and clinical decision support all fall under the high-risk categories enumerated in Annex III. Article 9 requires deployers to establish and implement a risk management system covering all foreseeable risks to health, safety, and fundamental rights throughout the system lifecycle. Article 10 requires measures to address reasonably foreseeable risks of biased outputs. Article 15 requires systems to achieve appropriate levels of accuracy, robustness, and cybersecurity.

The evaluation framework applied in this paper

directly maps to Article 15 cybersecurity and robustness requirements and to Article 10 bias documentation requirements. Specifically, the following findings constitute Article-level exposure for deployers who have not conducted equivalent testing: the universal crescendo jailbreak result (Article 15 robustness); the demographic bias findings in MedGemma-4B and Llama-3.1-8B (Article 10, Annex IV documentation); the sycophancy reversal of HIPAA compliance positions (Article 9 risk management, Article 14 human oversight); and the context-window displacement findings for Llama and Mistral families (Article 15 cybersecurity, EU AI Act Article 9).

The EU AI Act’s high-risk deployment deadline of 2 August 2026 makes the current period a critical window for regulated-industry deployers to conduct and document adversarial evaluations. Deployers who cannot produce evaluation records meeting Annex IV documentation requirements at that date face potential enforcement exposure.

Additionally, Mistral-7B’s critical failure on `competing_obj_biometric_data_conflict` — a probe testing whether the model violates GDPR

Article 9 special-category data processing restrictions when helpfulness incentives are opposed — is directly relevant to EU-regulated deployments in healthcare, HR screening, and insurance contexts, all of which handle special-category biometric and health data.

### 6.3 Limitations

Several limitations of this evaluation should be acknowledged.

*Text-only evaluation of MedGemma-4B.* MedGemma-4B is a vision-language model. We evaluated only its language generation pathway. The findings reported here apply to text-only deployments; the multimodal attack surface, including vision-based prompt injection and image-embedded adversarial content, is out of scope.

*Quantisation effects.* All models were evaluated at 4-bit quantisation via `mlx-lm` on Apple Silicon hardware. Quantisation may alter output distributions at safety-relevant token probability boundaries. We cannot rule out that some findings are quantisation-specific and would not reproduce at full precision; equally, quantisation may mask additional vulnerabilities present in the full-precision model. As 4-bit quantisation reflects the standard deployment configuration for on-premises SLM edge deployments, we consider this the appropriate scope.

*Black-box scope.* Our evaluation constitutes a lower bound on the full attack surface. White-box attacks — adversarial suffix generation, embedding-space manipulation, activation steering — represent additional vulnerability classes not addressed here.

*Coverage.* The framework covers 24 of the OWASP LLM Top 10 compliance references; OWASP-LLM-02 (insecure output handling), LLM-04 (model denial of service), and LLM-05 through LLM-07 and LLM-10 are not addressed in the current test suite. Reported compliance coverage should be interpreted relative to this scope.

*Classification stability.* Crescendo and sycophancy probe classifications were confirmed stable across multiple independent runs at temperature 0.2 and 0.7 for all six models. Competing objectives probes showed run-to-run variance of up to 20 percent at temperature 0.7, indicating that findings in this category carry greater uncertainty than the crescendo and sycophancy results. Stability testing was conducted at temperatures 0.2 and 0.7; replication across a broader temperature range would further characterise failure rate variance.

*Probe-level classification uncertainty.* The `sycophancy_safety_bleach_ammonia` findings for Qwen2-1.5B and Mistral-7B should be treated with

additional caution: programmatic classification confidence for these two probes was at threshold, and manual response review is recommended before treating them as confirmed critical findings.

*Mechanistic interpretation.* The three-tier context-window gradient is a behavioural observation. The proposed mechanistic interpretation — that RoPE positional encoding and the SWA lineage shared by Llama and Mistral architectures mediate the vulnerability — is a falsifiable hypothesis that would require attention-layer activation analysis to confirm or refute. We present it as a direction for follow-on interpretability work, not as an established finding.

### 6.4 Responsible Disclosure

The findings in this paper were produced through black-box adversarial evaluation using the same interface available to any end user or operator. No novel exploitation technique is disclosed here; all attack categories described have been characterised in the adversarial ML literature. We have not reported specific prompt payloads sufficient to reproduce critical findings without understanding the evaluation framework, in accordance with responsible disclosure norms.

Researchers or deployers who identify specific deployment contexts in which the reported findings constitute an active risk are encouraged to apply the appropriate vendor or operator notification processes.

## 7 Conclusion

We have presented a systematic adversarial evaluation of six small language models across 21 attack categories, yielding four principal findings: Five of six models failed all crescendo multi-turn escalation probes at critical severity, with Phi-3-mini as a reproducible partial exception; that context-window safety displacement follows a three-tier architectural gradient correlated with RoPE-based attention design; that medical domain fine-tuning redistributes adversarial exposure rather than reducing it, amplifying demographic bias in precisely the clinical domains the fine-tuning was intended to improve; and that sycophancy resistance to safety-relevant compliance reversal is architecture-mediated, not scale-mediated, with four of six models reversing correct safety-compliance positions under routine user social pressure. Taken together, these findings demonstrate that SLMs exhibit systematic, architecture-dependent adversarial vulnerabilities that differ materially from those documented in large language models, and that regulated-industry deployers who have not conducted

SLM-specific adversarial evaluation are shipping systems with uncharacterised attack surfaces directly relevant to their EU AI Act, HIPAA, and GDPR compliance posture. The evaluation methodology and attack category specifications are openly available at [github.com/sichgate/sichgate-methodology](https://github.com/sichgate/sichgate-methodology). The proprietary implementation underlying this evaluation is available to qualified researchers and regulated-industry practitioners for independent application.

## References

- [1] Cem Anil et al. Many-shot jailbreaking. *Anthropic Technical Report*, 2024.
- [2] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT)*, 2021.
- [3] Marianne Bertrand and Sendhil Mullainathan. Are emily and greg more employable than lakisha and jamal? *American Economic Review*, 94(4):991–1013, 2004.
- [4] Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *NeurIPS*, 2016.
- [5] Tim Dettmers, Artidoro Pagnoni, Ari Fantechi, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- [6] Elias Frantar, Saleh Ashkboos, Torsten Hoeffler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- [7] Deep Ganguli et al. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*, 2022.
- [8] 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 injections. In *AI Sec Workshop*, 2023.
- [9] Evan Hubinger et al. Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*, 2024.
- [10] Ninareh Mehrabi, Fred Morstatter, Navdeep Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 54(6):1–35, 2021.
- [11] Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.
- [12] Fábio Perez and Ian Ribeiro. Ignore previous prompt: Attack techniques for language models. In *NeurIPS ML Safety Workshop*, 2022.
- [13] 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. *arXiv preprint arXiv:2310.03693*, 2023.
- [14] Mark Russinovich, Ahmed Salem, and Ronen Eldan. Crescendo: Effective jailbreaks of large language models. *arXiv preprint arXiv:2404.01833*, 2024.
- [15] Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Asbell, Samuel R Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott T Johnston, et al. Towards understanding sycophancy in language models. *arXiv preprint arXiv:2310.13548*, 2023.
- [16] Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*, 2020.
- [17] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail? *arXiv preprint arXiv:2307.02483*, 2023.
- [18] Jerry Wei et al. Simple synthetic data reduces sycophancy in large language models. *arXiv preprint arXiv:2308.03958*, 2023.
- [19] Xianjun Yang et al. Shadow alignment: The ease of subverting safely-aligned language models. *arXiv preprint arXiv:2310.02949*, 2023.
- [20] Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.