

# **PANDAGUARD:** Systematic Evaluation of LLM Safety in the Era of Jailbreaking Attacks

### Anonymous Author(s) Affiliation Address email

## Abstract

1	Large language models (LLMs) have achieved remarkable capabilities but remain
2	vulnerable to adversarial prompts known as jailbreaks, which can bypass safety
3	alignment and elicit harmful outputs. Despite growing efforts in LLM safety
4	research, existing evaluations are often fragmented, focused on isolated attack or
5	defense techniques, and lack systematic, reproducible analysis. In this work, we
6	introduce PANDAGUARD, a unified and modular framework that models LLM
7	jailbreak safety as a multi-agent system comprising attackers, defenders, and judges.
8	Built on this framework, we develop PANDABENCH, a large-scale benchmark
9	encompassing over 50 LLMs, 20+ attack methods, 10+ defense mechanisms, and
0	multiple judgment strategies, requiring over 3 billion tokens to execute. Our
1	comprehensive evaluation reveals key insights into model vulnerabilities, defense
2	cost-performance trade-offs, and judge consistency. We find that no single defense
3	is optimal across all dimensions and that judge disagreement introduces nontrivial
4	variance in safety assessments. We release the full code, configurations, and
5	evaluation results to support transparent and reproducible research in LLM safety.

Homepage: https://panda-guard.github.io
 PANDAGUARD: https://github.com/Beijing-AISI/panda-guard
 PANDABENCH: https://hf.co/datasets/Beijing-AISI/panda-bench

## 18 1 Introduction

16

17

Large Language Models (LLMs), including architectures such as GPT, Llama, Qwen, and Gemini,
have achieved state-of-the-art performance across a wide range of natural language understanding
and generation tasks. Their rapid deployment in real-world applications—from content creation
and customer service to education and software development [1, 2]—highlights their transformative
potential. However, as LLMs become increasingly embedded in safety-critical systems, ensuring
their robustness and alignment has emerged as a paramount concern [3, 4, 5, 6, 7, 8].

A particularly acute threat to LLM safety is *jailbreaking*—a class of adversarial attacks in which
carefully engineered prompts circumvent alignment constraints and elicit harmful, biased, or unethical outputs [9, 10, 11, 12]. Successful jailbreaks can trigger toxic language, misinformation, or
even illegal instructions [13, 14, 15], undermining the guardrails built into state-of-the-art systems.
Accordingly, the development of robust defenses and rigorous evaluation protocols for LLM jailbreak
resistance has become an urgent research priority.

Despite valuable progress, current jailbreak evaluation approaches exhibit three key limitations. First, existing work often isolates individual components—such as attacks [10, 16, 17] or defenses [18,

19, 15, 20]—without capturing their systemic interplay. Second, there is a lack of standardized

benchmarking practices: evaluation protocols, datasets, and metrics remain fragmented [21, 22],

which hinders reproducibility and fair comparison. Third, most evaluations are conducted on a limited scale, covering only a small subset of models, threats, or response evaluators [23, 24]. Moreover,

critical factors such as computational cost, defense scalability, and the reliability of safety judges are

<sup>38</sup> often overlooked [25, 26].

Framework	#Attacks	#Defenses	#Judges	#Models	#LLM Interface
JAILJUDGE [27]	2	3	18	4	4 (HF, OpenAI, Gemini,)
EasyJailbreak [28]	12	0	7	10	4 (HF, OpenAI, kimi,)
AISafetyLab [21]	13	16	7	1	3 (HF,vLLM, OpenAI,)
HarmBench [29]	18	0	3	33	5 (HF, vLLM, OpenAI,)
SORRY-Bench [13]	0	0	1	56	1 (HF)
PANDAGUARD (Ours)	19	12	4	51	7 (HF, vLLM, SGLang, OpenAI,)

Table 1: Comparison of PANDAGUARD with existing LLM safety evaluation frameworks.

To address these challenges, we introduce PANDAGUARD, a unified and extensible evaluation framework that conceptualizes LLM jailbreak safety as a multi-agent system. In this formulation,

attackers, defenders, target models, and safety judges interact within a structured ecosystem, as 41 shwon in Figure 1. PANDAGUARD abstracts and modularizes each component, supporting plug-42 and-play experimentation with over 20 attack algorithms, 10+ defense mechanisms, and multiple 43 judgment strategies. This design facilitates controlled, reproducible evaluations and enables principled 44 analysis of cross-component trade-offs in model safety. Built atop PANDAGUARD, we further 45 develop PANDABENCH, a large-scale standardized benchmark suite encompassing approximately 3 46 billion tokens. PANDABENCH evaluates over 50 open and closed-source LLMs-spanning various 47 48 model sizes, release dates, and architectures—under diverse attack-defense combinations. Beyond breadth, our framework supports practical deployment via flexible user interaction modes (attack, 49 chat, serve) and compatibility with major inference engines including vLLM, SGLang, Ollama, 50 and remote APIs, thus enabling real-world usability and extensibility. Our contributions can be 51 summarized as follows: 52

- We propose PANDAGUARD, a principled multi-agent abstraction for LLM jailbreak safety
   that unifies attackers, defenders, target models, and judges within a modular system.
- We introduce PANDABENCH, a large-scale benchmark involving ~3B tokens and 50+ models, enabling broad and reproducible evaluations of jailbreak vulnerabilities and defenses.

• Through extensive empirical analysis, we uncover key insights into defense costeffectiveness, judge inconsistency, and evolving model vulnerabilities, offering actionable guidance for future safety research.

## 60 2 Background and Related Works

**Definitions.** Our framework conceptualizes jailbreaking as a multi-agent system with four distinct interacting components: attackers ( $\mathcal{A}$ ) generating adversarial prompts, target LLMs ( $\mathcal{M}$ ) processing inputs and generating outputs, defenders ( $\mathcal{D}$ ) implementing protection mechanisms, and safety judges ( $\mathcal{J}$ ) evaluating response safety. The primary objective of this system can be formalized as:

$$\min_{\mathcal{M},\mathcal{D}} \mathbb{E}_{P \sim \mathbf{P}, P' = \mathcal{A}(P)} [\mathcal{J}(\mathcal{D}(\mathcal{M}, P'))]$$
(1)

<sup>65</sup> Where *P* represents target prompts sampled from dataset **P**, *P'* is the adversarial prompt generated <sup>66</sup> by attacker  $\mathcal{A}$ ,  $\mathcal{M}$  is the target LLM, and  $\mathcal{D}$  represents defense mechanisms acting on either inputs or <sup>67</sup> outputs of  $\mathcal{M}$ . The safety judge  $\mathcal{J}$  typically outputs a binary value {0, 1} or a score in range [0, 1] <sup>68</sup> indicating whether a jailbreak was successful or its severity. While the overall objective involves <sup>69</sup> optimizing both models and defenses, our work primarily focuses on evaluating these components <sup>70</sup> within a unified framework. This formulation enables comprehensive analysis of safety dynamics, <sup>71</sup> emergent behaviors, and critical trade-offs between system components.

Jailbreak Attack Methodologies. Current attack methodologies can be categorized based on 72 their access level and strategic approach. From an access perspective, white-box attacks leverage 73 full knowledge of model parameters and architecture, utilizing gradient information to optimize 74 adversarial prompts, with GCG [10] pioneering gradient-based optimization of adversarial suffixes. 75 In contrast, black-box attacks operate without access to model parameters, exemplified by PAIR [17] 76 77 and AutoDAN [16]. Strategically, jailbreak attacks have evolved from static templates [11, 30, 31] to more adaptive approaches, including proxy-based optimization methods [32], zero-order 78 alternatives like random search [33] and genetic algorithms [16], and semantic-level attacks that 79 preserve malicious intent through linguistic transformations—such as PastTense [34], Rainbow 80 Teaming [35], ArtPrompt [36], and DeepInception [37], which uses nested fictional characters to 81 collectively work toward harmful goals, effectively bypassing safety mechanisms that focus primarily 82 on token-level patterns. 83

Defense Mechanisms. Defense mechanisms against jailbreak attacks span various implementation 84 approaches and processing strategies. System-level defenses operate externally to the model, imple-85 menting input filtering [38], response evaluation [39, 19], or in-context learning approaches [18, 40] 86 without requiring access to model parameters-including Self-Reminder [18], SmoothLLM [39] and 87 its semantic variant [41], perplexity filtering [38], paraphrasing techniques [38], and SelfDefend [42]. 88 In contrast, model-level defenses directly modify the LLM's parameters or internal representations, 89 through approaches like representation engineering [43, 15], adversarial training [29, 20], safety 90 fine-tuning [20, 44], RLHF [45], DPO [46], and Jailbreak Antidote [15]. An important but often over-91 looked aspect is the trade-off between security strength, computational overhead, and latency [15, 16], 92 93 which are critical for real-world deployment yet rarely evaluated systematically.

**Safety Evaluation.** Evaluating LLM safety presents significant challenges, particularly in establishing 94 consistent metrics and reliable judges. Current approaches include rule-based methods [10] that often 95 simply detect whether responses begin with refusals without necessarily assessing actual content 96 harmfulness. LLM-based judges [17] leverage other language models to evaluate response safety. 97 Human evaluation remains the gold standard, though it is highly resource-intensive and difficult 98 to scale. Recent studies have revealed concerning inconsistencies specifically with LLM-based 99 judges [47, 48], which can produce varying verdicts for identical inputs, raising questions about their 100 reliability as safety arbiters. These stability issues underscore the need for more robust methodologies 101 that can provide reliable assessments across diverse attack and defense scenarios. 102

Existing Benchmarks and Limitations. Several benchmarks have emerged to standardize jailbreak 103 evaluation, though each addresses only limited aspects of the safety ecosystem. JailbreakBench [24] 104 provides a centralized repository of adversarial prompts, HarmBench [29] implements various attacks 105 and defenses, SafetyBench [23] offers multiple-choice safety questions, and SORRY-Bench [13] 106 focuses on model refusal behaviors. EasyJailbreak [28] evaluates the effectiveness of various attack 107 methods against multiple models, and AISafetyLab [21] develops a tool for assessing model security 108 when both attack and defense methods are employed. JAILJUDGE [27] establishes a jailbreak 109 evaluation benchmark by integrating diverse attack scenarios and a multi-agent framework. Other 110 contributions include PromptBench [49], DecodingTrust [50], and TrustLLM [51], though they 111 primarily evaluate static templates rather than adaptive attacks. Despite these valuable efforts, 112 existing benchmarks suffer from key limitations: they often isolate specific attack vectors or defense 113 114 mechanisms rather than examining their interplay, lack standardized algorithm implementations that 115 lead to tight coupling between methods and models (obscuring true methodological contributions), 116 conduct evaluations at insufficient scale, and overlook critical aspects such as computational overhead and judge reliability. 117

Our work with PANDAGUARD addresses these limitations by providing a comprehensive framework that integrates the full spectrum of components in the jailbreaking ecosystem within a standardized evaluation protocol. By enabling systematic variation of models, attacks, defenses, and evaluation methods, PANDAGUARD facilitates rigorous, reproducible research. Building on this foundation, PANDABENCH implements extensive evaluations at scale to support the development of robust safety mechanisms that balance security, efficiency, and user experience.

## 124 3 PANDAGUARD: A Framework for Safety-Critical LLM Evaluation

In this section, we introduce PANDAGUARD, a modular and extensible framework designed to address
key limitations in existing LLM jailbreak safety evaluations. While prior efforts often focus on
isolated attacks or defenses, PANDAGUARD systematically models the full safety pipeline—including
attackers, defenders, target models, and judges—within a unified, reproducible environment. The
framework enables researchers to study the intricate interactions and trade-offs between components,
fostering a deeper understanding of safety dynamics across diverse settings.

Architecture. PANDAGUARD uses a pipeline-based design to orchestrate interactions among system components, as formalized in Equation 1. Upon receiving a target prompt (e.g., a jailbreak goal), the system invokes configurable attack modules to generate adversarial queries. These queries are processed by defense mechanisms, which may modify the input or filter the output before reaching the target LLM. The generated responses are then assessed by one or more safety judges to determine whether harmful content was successfully elicited.



Figure 1: The PANDAGUARD framework architecture illustrating the end-to-end pipeline for LLM safety evaluation. The system connects three key components: Attackers, Defenders, and Judges. The framework supports diverse LLM interfaces and demonstrates several practical applications including interactive chat, API serving, attack generation, and systematic evaluation.

This modular architecture enables controlled experimentation by allowing researchers to fix any
component and systematically vary others. For example, one can evaluate a defense strategy across
multiple attacks, compare LLM vulnerabilities under a common threat model, or deploy defenseenhanced LLMs in interactive settings. The use of standardized interfaces across all components
ensures both scalability and reproducibility.

**Component Abstraction and Implementation.** PANDAGUARD provides consistent abstraction layers across all modules. For attackers, we define a base interface with an attack() method that transforms user queries into adversarial prompts. Our implementation supports a wide range of methods, including black-box attacks (e.g., PAIR [17], DeepInception [37], AutoDAN [16]) and optimization-based techniques such as GCG [10] and RandomSearch [33].

Defender modules implement a defense() method to process potentially harmful content. We
support three major paradigms: (1) detection-based methods (e.g., PerplexityFilter [38]) that filter
adversarial prompts, (2) prompt-based defenses (e.g., SelfReminder [18], GoalPriority [40], Smooth-

LLM [39]) that manipulate input phrasing, and (3) representation-level methods such as Jailbreak

151 Antidote [15] that adjust internal model states to neutralize threats.

The target LLM interface supports both commercial API-based models (e.g., OpenAI, Anthropic, Gemini) and locally hosted models via frameworks like vLLM [52], SGLang [53], Ollama [54], and Transformers [55]. Key functionalities include generate(), evaluate\_log\_likelihood(), and

<sup>155</sup> batch\_generate(), ensuring consistent behavior across backends.

156 Safety judges implement a judge() method to evaluate responses using standardized scoring proto-157 cols. PANDAGUARD supports both rule-based judges [10] and LLM-based judges [17, 56], enabling

<sup>158</sup> comparative analysis of judgment consistency and reliability.

**Configuration-Driven Experimentation.** PANDAGUARD uses YAML-based configuration files to specify pipeline components, hyperparameters, and evaluation options without requiring code

changes. This design facilitates reproducible experimentation and transparent system specification.

162 Code 3.1 and Code 3.2 illustrate typical usage patterns.



<sup>163</sup> This configuration system supports precise and scalable experimentation. New components can be <sup>164</sup> registered via entry points, enabling extensibility without modifying core logic.

Versatile Interface Options. PANDAGUARD offers multiple modes for research and deployment.
 The command-line interface supports commands such as panda-guard chat, serve, inference, and attack, enabling users to deploy defense-enhanced LLMs, run interactive sessions, launch
 API services, or conduct targeted jailbreak generation. The design is optimized for integration into
 real-world research pipelines and production environments.

PANDAGUARD serves as the technical foundation of PANDABENCH, enabling the most compre hensive integration of jailbreak attacks, defenses, and evaluators to date. Its extensible architecture,
 multi-backend support, and reproducibility features make it a powerful tool for both academic
 research and practical LLM safety evaluation.

## **4 PANDABENCH: Empirical Results and Key Insights**

To comprehensively evaluate LLM jailbreak safety, we build PANDABENCH atop the PANDAGUARD
framework. Unlike previous benchmarks that focus on limited models, isolated attack methods,
or omit defense and judge considerations [24, 26], PANDABENCH offers the most comprehensive
and reproducible jailbreak safety evaluation to date. PANDABENCH includes 51 diverse LLMs
across model families and scales, 18 attack algorithms, 9 defense mechanisms, and multiple judging
strategies. We adopt Attack Success Rate (ASR) as the primary metric, following the PAIR [17]
criterion—an attack is deemed successful only if the judge assigns a maximum score of 10.

To ensure fair comparison, we unify the proxy model used in attack and defense interactions. Specifically, we use Llama-3.1-8B [57] to generate adversarial prompts for attack algorithms and act as the agent for defense mechanisms. This eliminates discrepancies introduced by different backbone models and ensures that observed performance differences arise solely from algorithmic design. The full experimental setup, including all scripts, configurations, and examples, is available in our HuggingFace repository. https://hf.co/datasets/Beijing-AISI/panda-bench.



Figure 2: **Model-wise safety analysis.** (a) ASR vs. release date for various LLMs. (b) ASR across different harm categories with and without defense mechanisms. (c) Overall ASR for all evaluated LLMs with and without defense mechanisms.

### 188 4.1 Model-wise Safety Analysis

Figure 2 illustrates how various LLMs respond to jailbreaking attempts, revealing vulnerability
 patterns both with and without defensive countermeasures in place. The visualization captures model
 safety across multiple dimensions.

**Safety trends across model evolution.** Figure 2a plots ASR versus release date for a range of LLMs, 192 revealing multiple important patterns. First, we observe substantial variation in safety performance 193 across model families, with proprietary models like GPT and Claude generally exhibiting lower ASRs 194 compared to open-source models, reflecting stronger safety alignment. However, safety does not 195 196 consistently improve over time-in fact, the variance in safety performance increases in newer models. This indicates that improvements in safety do not necessarily align with general model capabilities 197 198 but are likely influenced by specific alignment strategies used during development. Additionally, within the same generation, larger models tend to exhibit better safety properties, but newer models 199 (e.g., Owen3) can have worse safety performance than older versions (e.g., Owen2.5), highlighting 200 that safety is not an emergent property of scale or recency but requires deliberate optimization. 201

Vulnerability across harm categories. Figure 2b breaks down ASR by harm category, comparing performance with and without defenses. While defense mechanisms reduce vulnerability in all categories, some harm categories (e.g., malware/hacking, fraud/deception, privacy) remain more difficult to mitigate, even with defenses in place. This suggests that certain types of harm may be underrepresented in alignment training data or may be inherently more difficult to defend against due to the need to provide benign yet related information (e.g., explaining cybersecurity concepts without enabling malicious activities).

**Defense impact across models.** Figure 2c summarizes overall ASRs for all evaluated models, both with and without defense mechanisms. We observe that defenses consistently reduce ASRs by approximately one-third to one-half, with more significant gains for models with higher initial



Figure 3: ASR heatmap for different attack methods against various LLMs. Higher values indicate more successful attacks.

vulnerability. For example, Claude-3.5 and GPT-40 exhibit ASRs below 3% without defenses,
while DeepSeek-R1 and Qwen3-1.7B exceed 20% without defenses. This highlights gaps in safety
alignment strategies across different model providers and emphasizes the importance of both inherent
model safety and additional defense mechanisms.

### 216 4.2 Attack and Defense Mechanisms Analysis

217 We now turn to a key aspect of LLM safety: the interplay between attack and defense mechanisms,

and the inherent trade-offs involved in designing robust defenses. Our analysis, visualized in

Figure 4, spans three dimensions: attack-defense effectiveness (a), computational overhead (b), and

220 performance impact (c).



Figure 4: Attack and defense mechanisms analysis. (a) Heatmap of attack success rates across different combinations of attack and defense methods. (b) Trade-off between defense effectiveness and computational overhead measured in total tokens. (c) Trade-off between defense effectiveness and impact on model performance as measured by Alpaca winrate [58].

Attack-defense interaction patterns. Figure 4a presents a comprehensive heatmap of ASRs for each attack-defense pair, with both axes sorted by average effectiveness. All defense methods consistently reduce ASR compared to the Baseline (no defense), reaffirming their necessity in safetycritical settings. Semantic-level defenses, such as Paraphrase [38] and Semantic SmoothLLM [41], demonstrate exceptional robustness, maintaining low ASRs even against sophisticated attacks by semantically rewriting prompts and uncovering user intent.

Among attack strategies, ensemble methods like RandomSearch [33] perform most effectively, maintaining high ASRs despite countermeasures. Intent-masking techniques such as PastTense [34] and AutoDAN [16] also prove potent by disguising harmful intent. Notably, even template-based approaches like AIM and BETTER\_DAN [30] exhibit measurable success, revealing persistent blind spots in existing defenses.

**Computational efficiency of defenses.** Figure 4b explores the trade-off between defense effectiveness and computational cost. Given the diversity in helper models and hardware setups, we standardize evaluations by measuring token overhead relative to a common Baseline. Dialog-based defenses such as SmoothLLM [39] and Semantic SmoothLLM [41] incur significantly higher token usage—up to 5x the Baseline—due to multi-turn interactions.

This raises practical concerns for deployment. The green reference line represents the estimated token cost of post-generation safety filtering using PAIR's [17] judge prompt, providing a cost-efficiency benchmark. These results underscore the need to balance defense strength with operational feasibility.

Performance impact of defenses. Figure 4c evaluates the effect of defenses on standard model utility using AlpacaEval winrate [58]. Stronger defenses tend to reduce performance more severely, with some methods degrading output quality by up to 25%. This highlights a critical challenge: preserving model usefulness while ensuring safety.

Taken together, these findings reveal a three-way trade-off among safety effectiveness, computational
efficiency, and task performance. For instance, SelfDefense [42] excels in reducing ASR but causes
notable performance drops, whereas PerplexityFilter [38] retains high utility but offers weaker
protection. Paraphrase [38] offers a practical middle ground with balanced trade-offs.

Overall, this section emphasizes the importance of holistic evaluation frameworks like PANDAGUARD,
 which support systematic, reproducible comparisons across safety mechanisms, guiding the design of
 well-balanced defenses for real-world deployment.

#### 251 4.3 Safety Judge Reliability Analysis

Since jailbreak detection hinges critically on judge reliability, inconsistencies among evaluation
 strategies may distort our understanding of model vulnerabilities and defense effectiveness. Our
 final analysis examines the reliability and consistency of different safety judgment methodologies.
 Figure 5 presents our findings on judge reliability and agreement.

(a) Social Adult ICL GoalPriority AntPrompt Cipher GPT (b) - 1.000 0.0



Figure 5: **Safety judge reliability analysis.** (a) Radar charts comparing ASR judgments by different judges across harm categories, defense methods, and attack methods. Judges include rule-based and LLM-based (GPT-40, Qwen2.5, Llama3.3). (b) Cohen's Kappa matrix showing agreement between different judges.

Judge behavior variability. Figure 5a reveals substantial variations in how different judges evaluate the same model outputs. The rule-based judge (GCG [10]), which primarily detects refusal patterns, consistently reports significantly higher ASRs across all harm categories, defense methods, and attack strategies compared to LLM-based judges. This discrepancy underscores differing philosophies: whether jailbreaks are defined by refusal absence or by actual content harm.

Among LLM-based judges, we observe coherent evaluation patterns with instructive variations in sensitivity. GPT-40 demonstrates greater stringency in safety evaluations compared to Qwen2.5 and Llama3.3, particularly for categories like expert advice and government decision-making. This
 diversity, even when using identical judging prompts, provides valuable perspectives that enrich our
 understanding of safety boundaries across model families.

The radar charts further reveal category-specific judge behaviors. For instance, all judges show relatively higher agreement on sexual/adult content and harassment/discrimination, while exhibiting greater divergence on categories like malware/hacking and economic harm. This pattern may reflect varying levels of clarity in safety guidelines across different harm types, as well as differing interpretations of what constitutes harmful information versus legitimate educational content.

**Inter-judge agreement analysis.** Figure 5b quantifies the relationships between different judges using Cohen's Kappa coefficients. The distinct approach of the rule-based judge (GCG) compared to LLM-based judges is reflected in Kappa values ranging from 0.071 to 0.126.

The moderate agreement among LLM-based judges reveals significant challenges in safety evaluation. While there is some consensus on extreme cases, boundary judgments vary considerably, reflecting the inherent difficulty in defining harmful content across different contexts, cultures, and use cases. The varying sensitivities—even between sophisticated models like GPT-40 and Qwen2.5—highlight the subjective nature of harm assessment and the absence of universal standards.

These findings underscore the complexity of safety evaluation and the limitations of relying on single judgment sources. What one system deems harmful might be considered educational or contextually appropriate by another. This variability emphasizes the need for frameworks like PANDAGUARD that support multi-dimensional assessment approaches. By enabling controlled comparison of different judging strategies, our benchmark provides researchers with insights into evaluation reliability and helps advance more nuanced, context-aware safety assessment methodologies that acknowledge these fundamental challenges.

## 286 5 Conclusion

This work presents PANDAGUARD, a unified and extensible framework for systematically evaluating 287 the jailbreak robustness of large language models. By modeling the safety ecosystem as a multi-288 agent interaction among attackers, defenders, target models, and judges, PANDAGUARD enables 289 modular experimentation, reproducibility, and in-depth analysis across the full spectrum of safety 290 components. Built on this framework, PANDABENCH conducts the most comprehensive empirical 291 study to date, spanning over 50 LLMs, 20 attack techniques, and 10 defense strategies. Our findings 292 reveal nuanced trade-offs among safety, cost, and performance; expose reliability challenges in 293 current safety judgments; and offer actionable insights for the design of more balanced and effective 294 safety mechanisms. We release the full framework, benchmark suite, and evaluation results to foster 295 transparent, reproducible, and forward-looking research in LLM safety. 296

## <sup>297</sup> A Technical Appendices and Supplementary Material

Technical appendices with additional results, figures, graphs and proofs may be submitted with the paper submission before the full submission deadline (see above), or as a separate PDF in the ZIP file below before the supplementary material deadline. There is no page limit for the technical appendices.

## 302 **References**

- [1] Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh
   Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward
   design via coding large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- [2] Oguzhan Topsakal and Tahir Cetin Akinci. Creating large language model applications utilizing
   langchain: A primer on developing llm apps fast. In *International Conference on Applied Engineering and Natural Sciences*, volume 1, pages 1050–1056, 2023.
- [3] Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor,
   Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. Taxonomy of risks posed

- by language models. In *Proceedings of the 2022 ACM conference on fairness, accountability, and transparency*, pages 214–229, 2022.
- [4] Haiyang Wang, Yihao Li, Yue Wang, Pan Liu, and Pengxiao Li. Navigating the risks: A review
   of safety issues in large language models. In 2024 IEEE 24th International Conference on
   Software Quality, Reliability, and Security Companion (QRS-C), pages 74–83. IEEE, 2024.
- [5] Yifan Yao, Jinhao Duan, Kaidi Xu, Yuanfang Cai, Zhibo Sun, and Yue Zhang. A survey on large
   language model (llm) security and privacy: The good, the bad, and the ugly. *High-Confidence Computing*, page 100211, 2024.
- [6] Kaiwen Zhou, Chengzhi Liu, Xuandong Zhao, Anderson Compalas, and Xin Eric Wang.
   Multimodal situational safety. In Workshop on Responsibly Building the Next Generation of Multimodal Foundational Models, 2024.
- [7] Yi Zeng, Yu Yang, Andy Zhou, Jeffrey Ziwei Tan, Yuheng Tu, Yifan Mai, Kevin Klyman,
   Minzhou Pan, Ruoxi Jia, Dawn Song, et al. Air-bench 2024: A safety benchmark based on
   regulation and policies specified risk categories. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [8] Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and
   Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. In *The Twelfth International Conference on Learning Representations*, 2024.
- [9] Junjie Chu, Yugeng Liu, Ziqing Yang, Xinyue Shen, Michael Backes, and Yang Zhang. Com prehensive assessment of jailbreak attacks against llms. *arXiv preprint arXiv:2402.05668*, 2024.
- [10] Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and transferable
   adversarial attacks on aligned language models, 2023.
- [11] Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei
   Zhang, Kailong Wang, and Yang Liu. Jailbreaking chatgpt via prompt engineering: An empirical
   study. *arXiv preprint arXiv:2305.13860*, 2023.
- [12] Sibo Yi, Yule Liu, Zhen Sun, Tianshuo Cong, Xinlei He, Jiaxing Song, Ke Xu, and Qi Li.
   Jailbreak attacks and defenses against large language models: A survey. *arXiv preprint arXiv:2407.04295*, 2024.
- [13] Tinghao Xie, Xiangyu Qi, Yi Zeng, Yangsibo Huang, Udari Madhushani Sehwag, Kaixuan
   Huang, Luxi He, Boyi Wei, Dacheng Li, Ying Sheng, et al. Sorry-bench: Systematically
   evaluating large language model safety refusal behaviors. *arXiv preprint arXiv:2406.14598*,
   2024.
- [14] Youliang Yuan, Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Jiahao Xu, Tian Liang, Pinjia
   He, and Zhaopeng Tu. Refuse whenever you feel unsafe: Improving safety in Ilms via decoupled
   refusal training. *arXiv preprint arXiv:2407.09121*, 2024.
- [15] Guobin Shen, Dongcheng Zhao, Yiting Dong, Xiang He, and Yi Zeng. Jailbreak antidote:
   Runtime safety-utility balance via sparse representation adjustment in large language models.
   *The Thirteenth International Conference on Learning Representations*, 2025.
- [16] Xiaogeng Liu, Nan Xu, Muhao Chen, and Chaowei Xiao. Autodan: Generating stealthy
   jailbreak prompts on aligned large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric
   Wong. Jailbreaking black box large language models in twenty queries. In *R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*, 2023.
- [18] Yueqi Xie, Jingwei Yi, Jiawei Shao, Justin Curl, Lingjuan Lyu, Qifeng Chen, Xing Xie, and
   Fangzhao Wu. Defending chatgpt against jailbreak attack via self-reminders. *Nature Machine Intelligence*, 5(12):1486–1496, 2023.

- [19] Yihan Wang, Zhouxing Shi, Andrew Bai, and Cho-Jui Hsieh. Defending llms against jailbreak ing attacks via backtranslation. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 16031–16046, 2024.
- [20] Jiaming Ji, Donghai Hong, Borong Zhang, Boyuan Chen, Josef Dai, Boren Zheng, Tianyi Qiu,
   Boxun Li, and Yaodong Yang. Pku-saferlhf: A safety alignment preference dataset for llama
   family models. *arXiv e-prints*, pages arXiv–2406, 2024.
- [21] Zhexin Zhang, Leqi Lei, Junxiao Yang, Xijie Huang, Yida Lu, Shiyao Cui, Renmiao Chen,
   Qinglin Zhang, Xinyuan Wang, Hao Wang, et al. Aisafetylab: A comprehensive framework for
   ai safety evaluation and improvement. *arXiv preprint arXiv:2502.16776*, 2025.
- [22] Zhao Xu, Fan Liu, and Hao Liu. Bag of tricks: Benchmarking of jailbreak attacks on llms. In
   *Advances in Neural Information Processing Systems*, 2024.
- [23] Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu,
   Xuanyu Lei, Jie Tang, and Minlie Huang. Safetybench: Evaluating the safety of large language
   models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15537–15553, 2024.
- Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco
   Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian
   Tramèr, Hamed Hassani, and Eric Wong. Jailbreakbench: An open robustness benchmark for
   jailbreaking large language models. In *NeurIPS Datasets and Benchmarks Track*, 2024.
- [25] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen,
   Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language
   models. ACM transactions on intelligent systems and technology, 15(3):1–45, 2024.
- <sup>382</sup> [26] Davide Biarese. Advbench: a framework to evaluate adversarial attacks against fraud detection <sup>383</sup> systems. 2022.
- [27] Fan Liu, Yue Feng, Zhao Xu, Lixin Su, Xinyu Ma, Dawei Yin, and Hao Liu. Jailjudge: A
   comprehensive jailbreak judge benchmark with multi-agent enhanced explanation evaluation
   framework, 2024.
- [28] Weikang Zhou, Xiao Wang, Limao Xiong, Han Xia, Yingshuang Gu, Mingxu Chai, Fukang
   Zhu, Caishuang Huang, Shihan Dou, Zhiheng Xi, et al. Easyjailbreak: A unified framework for
   jailbreaking large language models. *arXiv preprint arXiv:2403.12171*, 2024.
- [29] Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham
   Sakhaee, Nathaniel Li, Steven Basart, Bo Li, et al. Harmbench: a standardized evaluation frame work for automated red teaming and robust refusal. In *Proceedings of the 41st International Conference on Machine Learning*, pages 35181–35224, 2024.
- [30] Alex Albert. Jailbreak chat. https://jailbreakchat-hko42cs2r-alexalbertt-s-team.
   vercel.app/, 2025. Accessed: 2025-05-11.
- [31] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does LLM safety
   training fail? In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [32] Qizhang Li, Yiwen Guo, Wangmeng Zuo, and Hao Chen. Improved generation of adversarial
   examples against safety-aligned LLMs. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [33] Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Jailbreaking leading
   safety-aligned llms with simple adaptive attacks. In *ICML 2024 Next Generation of AI Safety Workshop*, 2024.
- [34] Maksym Andriushchenko and Nicolas Flammarion. Does refusal training in llms generalize to
   the past tense? In *Neurips Safe Generative AI Workshop 2024*, 2024.

- [35] Mikayel Samvelyan, Sharath Chandra Raparthy, Andrei Lupu, Eric Hambro, Aram Markosyan,
   Manish Bhatt, Yuning Mao, Minqi Jiang, Jack Parker-Holder, Jakob Foerster, et al. Rain bow teaming: Open-ended generation of diverse adversarial prompts. *Advances in Neural Information Processing Systems*, 37:69747–69786, 2024.
- [36] Fengqing Jiang, Zhangchen Xu, Luyao Niu, Zhen Xiang, Bhaskar Ramasubramanian, Bo Li,
   and Radha Poovendran. Artprompt: Ascii art-based jailbreak attacks against aligned llms. In
   *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 15157–15173, 2024.
- [37] Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. Deepinception:
   Hypnotize large language model to be jailbreaker. In *Neurips Safe Generative AI Workshop* 2024, 2024.
- [38] Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping-yeh
   Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. Baseline de fenses for adversarial attacks against aligned language models. *arXiv preprint arXiv:2309.00614*, 2023.
- [39] Alexander Robey, Eric Wong, Hamed Hassani, and George J Pappas. Smoothllm: Defending
   large language models against jailbreaking attacks. *arXiv preprint arXiv:2310.03684*, 2023.
- [40] Zhexin Zhang, Junxiao Yang, Pei Ke, Fei Mi, Hongning Wang, and Minlie Huang. Defending
   large language models against jailbreaking attacks through goal prioritization. In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long
   Papers), pages 8865–8887, 2024.
- [41] Jiabao Ji, Bairu Hou, Alexander Robey, George J Pappas, Hamed Hassani, Yang Zhang, Eric
   Wong, and Shiyu Chang. Defending large language models against jailbreak attacks via semantic
   smoothing. arXiv preprint arXiv:2402.16192, 2024.
- [42] Mansi Phute, Alec Helbling, Matthew Daniel Hull, ShengYun Peng, Sebastian Szyller, Cory
   Cornelius, and Duen Horng Chau. Llm self defense: By self examination, llms know they are
   being tricked. In *The Second Tiny Papers Track at ICLR 2024*, 2024.
- [43] Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,
   Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A
   top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.
- [44] Chia-Yi Hsu, Yu-Lin Tsai, Chih-Hsun Lin, Pin-Yu Chen, Chia-Mu Yu, and Chun-Ying Huang.
   Safe lora: The silver lining of reducing safety risks when finetuning large language models.
   Advances in Neural Information Processing Systems, 37:65072–65094, 2024.
- [45] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
   Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton,
   Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano,
   Jan Leike, and Ryan Lowe. Training language models to follow instructions with human
   feedback. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors,
   Advances in Neural Information Processing Systems, 2022.
- [46] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and
   Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
   *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [47] Tianhao Wu, Weizhe Yuan, Olga Golovneva, Jing Xu, Yuandong Tian, Jiantao Jiao, Jason
   Weston, and Sainbayar Sukhbaatar. Meta-rewarding language models: Self-improving alignment
   with llm-as-a-meta-judge. *arXiv preprint arXiv:2407.19594*, 2024.
- [48] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li,
   Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*, 2024.

- [49] Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi
  Yang, Wei Ye, Yue Zhang, Neil Zhenqiang Gong, et al. Promptbench: Towards evaluating the
  robustness of large language models on adversarial prompts. *arXiv e-prints*, pages arXiv–2306,
  2023.
- [50] Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian
   Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment
   of trustworthiness in gpt models. In *NeurIPS*, 2023.
- [51] Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang,
   Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. Trustllm: Trustworthiness in large language models.
   *arXiv preprint arXiv:2401.05561*, 3, 2024.
- [52] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu,
   Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large lan guage model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium* on Operating Systems Principles, 2023.
- [53] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Livia Sun, Jeff Huang, Cody Hao Yu,
   Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. Sglang: Efficient execution
   of structured language model programs. *Advances in Neural Information Processing Systems*,
   37:62557–62583, 2024.
- 472 [54] Ollama Team. Ollama. https://ollama.com, 2023. Accessed: 2025-05-08.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony
  Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer,
  Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain
  Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-theart natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020.
  Association for Computational Linguistics.
- [56] Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron
   Singer, and Amin Karbasi. Tree of attacks: Jailbreaking black-box llms automatically. *Advances in Neural Information Processing Systems*, 37:61065–61105, 2024.
- [57] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian,
  Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama
  3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [58] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy
   Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following
   models. https://github.com/tatsu-lab/alpaca\_eval, 5 2023.

## **189** NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

**The checklist answers are an integral part of your paper submission.** They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. 505 While "[Yes] " is generally preferable to "[No] ", it is perfectly acceptable to answer "[No] " provided a 506 proper justification is given (e.g., "error bars are not reported because it would be too computationally 507 expensive" or "we were unable to find the license for the dataset we used"). In general, answering 508 "[No] " or "[NA] " is not grounds for rejection. While the questions are phrased in a binary way, we 509 acknowledge that the true answer is often more nuanced, so please just use your best judgment and 510 write a justification to elaborate. All supporting evidence can appear either in the main paper or the 511 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification 512 please point to the section(s) where related material for the question can be found. 513

- 514 IMPORTANT, please:
- Delete this instruction block, but keep the section heading "NeurIPS Paper Checklist",
- Keep the checklist subsection headings, questions/answers and guidelines below.
- Do not modify the questions and only use the provided macros for your answers.
- 518 1. Claims
- 519 Question: Do the main claims made in the abstract and introduction accurately reflect the 520 paper's contributions and scope?
- 521 Answer: **[TODO]**
- 522 Justification: [TODO]
  - Guidelines:

523

524

525

526

527

528

529

530

531

532

533

534

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
  - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
  - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

- 535 Answer: [TODO]
- 536 Justification: **[TODO]**

537	Guidelines:
538 539	• The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
540	• The authors are encouraged to create a separate "Limitations" section in their paper.
541	• The paper should point out any strong assumptions and how robust the results are to
542	violations of these assumptions (e.g., independence assumptions, noiseless settings,
543	model well-specification, asymptotic approximations only holding locally). The authors
544	should reflect on how these assumptions might be violated in practice and what the
545	implications would be.
546	• The authors should reflect on the scope of the claims made, e.g., if the approach was
547 548	only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
549	• The authors should reflect on the factors that influence the performance of the approach.
550	For example, a facial recognition algorithm may perform poorly when image resolution
551	is low or images are taken in low lighting. Or a speech-to-text system might not be
552	used reliably to provide closed captions for online lectures because it fails to handle
553	technical jargon.
554 555	• The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
556 557	• If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
558	• While the authors might fear that complete honesty about limitations might be used by
559	reviewers as grounds for rejection, a worse outcome might be that reviewers discover
560	limitations that aren't acknowledged in the paper. The authors should use their best
561	judgment and recognize that individual actions in favor of transparency play an impor-
562	will be specifically instructed to not penalize honesty concerning limitations
564 <b>3</b> .	<ul> <li>Theory assumptions and proofs</li> </ul>
565	Ouestion: For each theoretical result, does the paper provide the full set of assumptions and
566	a complete (and correct) proof?
567	Answer: [TODO]
568	Justification: [TODO]
569	Guidelines:
570	• The answer NA means that the paper does not include theoretical results.
571	• All the theorems, formulas, and proofs in the paper should be numbered and cross-
572	referenced.
573	• All assumptions should be clearly stated or referenced in the statement of any theorems.
574	• The proofs can either appear in the main paper or the supplemental material, but if
575	they appear in the supplemental material, the authors are encouraged to provide a short
576	proof sketch to provide intuition.
577	• Inversely, any informal proof provided in the core of the paper should be complemented
578	by formal proofs provided in appendix or supplemental material.
579	• Theorems and Lemmas that the proof relies upon should be properly referenced.
580 4.	Experimental result reproducibility
581	Question: Does the paper fully disclose all the information needed to reproduce the main ex-
582	of the paper (regardless of whether the code and data are provided or not)?
503	or the paper (regardless of whether the code and data are provided or not):
584	Answer: [TODO]
585	Justification: [TODO]
500	
586	Guidelines:

588 589 590	• If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
591 592	• If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
502	• Depending on the contribution reproducibility can be accomplished in various ways
593 594	For example, if the contribution is a novel architecture, describing the architecture fully
595	might suffice, or if the contribution is a specific model and empirical evaluation, it may
596	be necessary to either make it possible for others to replicate the model with the same
597	dataset, or provide access to the model. In general, releasing code and data is often
598	one good way to accomplish this, but reproducibility can also be provided via detailed
599	instructions for how to replicate the results, access to a hosted model (e.g., in the case
600 601	of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
602	• While NeurIPS does not require releasing code, the conference does require all submis-
603	sions to provide some reasonable avenue for reproducibility, which may depend on the
604	nature of the contribution. For example
605	(a) If the contribution is primarily a new algorithm, the paper should make it clear how
606	
607 608	(b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
609	(c) If the contribution is a new model (e.g., a large language model), then there should
610	either be a way to access this model for reproducing the results or a way to reproduce
611	the model (e.g., with an open-source dataset or instructions for how to construct
612	the dataset).
613	(d) We recognize that reproducibility may be tricky in some cases, in which case
614	authors are welcome to describe the particular way they provide for reproducibility.
615	In the case of closed-source models, it may be that access to the model is limited in
616	some way (e.g., to registered users) but it should be possible for other researchers
	some may (e.g., to registered users), our it should be possible for other researchers
617	to have some path to reproducing or verifying the results.
617 618	<ul><li>to have some path to reproducing or verifying the results.</li><li>5. Open access to data and code</li></ul>
617 618 619	<ul> <li>to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instruc-</li></ul>
617 618 619 620	<ul> <li>5. Open access to data and code</li> <li>Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental</li> </ul>
617 618 619 620 621	<ul> <li>5. Open access to data and code</li> <li>Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?</li> </ul>
617 618 619 620 621 622	<ul> <li>5. Open access to data and code</li> <li>Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?</li> <li>Answer: [TODO]</li> </ul>
617 618 619 620 621 622 623	<ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO]</li></ul>
617 618 619 620 621 622 623 624	<ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines:</li></ul>
617 618 619 620 621 622 623 624 625	<ul> <li>5. Open access to data and code</li> <li>Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?</li> <li>Answer: [TODO]</li> <li>Justification: [TODO]</li> <li>Guidelines:</li> <li>The answer NA means that paper does not include experiments requiring code.</li> </ul>
617 618 619 620 621 622 623 624 625 626	<ul> <li>bonc may (e.g., to registered decis), our national of possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code <ul> <li>Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?</li> <li>Answer: [TODO]</li> <li>Justification: [TODO]</li> <li>Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/</li> </ul> </li> </ul></li></ul>
617 618 619 620 621 622 623 623 624 625 626 627	<ul> <li>bound may (eig., to regulated users), our national of possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> </ul></li></ul>
617 618 619 620 621 622 623 623 624 625 626 627 628	<ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be</li> </ul></li></ul>
617 618 619 620 621 622 623 624 625 625 626 627 628 629	<ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not</li> </ul></li></ul>
617 618 619 620 621 622 623 624 625 626 627 628 629 630	<ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source </li> </ul></li></ul>
617 618 620 621 622 623 624 625 626 625 626 627 628 629 630 631	<ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> </ul></li></ul>
617 618 620 621 622 623 624 625 626 627 628 629 630 631 632	<ul> <li>bound may (e.g., to regulated users), our national de possible for outer researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). <ul> <li>The instructions should contain the exact command and environment needed to run to</li> </ul></li></ul></li></ul>
617 618 620 621 622 623 623 624 625 626 627 628 629 630 631 632 633	<ul> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/benchmark).</li> </ul></li></ul>
617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634	<ul> <li>bond and (e.g., or registered does), our restored or possible for other researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).</li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> </ul></li></ul>
617 618 619 620 621 622 623 623 624 625 626 627 628 629 630 631 632 633 634 635	<ul> <li>bond may (e.g., to regulatered abela), our nanodative possible for outer researchers to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).</li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> </ul></li></ul>
617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636	<ul> <li>to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li></ul></li></ul>
617 618 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637	<ul> <li>to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc. The authors should provide scripts to reproduce all experimental results for the new</li></ul></li></ul>
617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638	<ul> <li>to have some path to reproducing or verifying the results.</li> <li>5. Open access to data and code Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). <ul> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. <ul> <li>The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li> <li>The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they </li> </ul></li></ul></li></ul></li></ul>
617         618         619         620         621         622         623         624         625         626         627         628         630         631         632         633         634         635         636         637         638         639	<ul> <li>to have some path to reproducing or verifying the results.</li> <li><b>Open access to data and code</b> Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). </li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details. The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc. The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.</li></ul></li></ul>
617         618         619         620         621         622         623         624         625         626         627         628         629         630         631         632         633         634         635         636         637         638         639         640	<ul> <li>base and (e.g., to reproducing or verifying the results.</li> <li><b>5. Open access to data and code</b> Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material? Answer: [TODO] Justification: [TODO] Guidelines: <ul> <li>The answer NA means that paper does not include experiments requiring code.</li> <li>Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark). <ul> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.</li> <li>The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.</li> <li>The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why. <ul> <li>At submission time, to preserve anonymity, the authors should release anonymized</li> </ul></li></ul></li></ul></li></ul>

642 643		• Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.
644	6.	Experimental setting/details
645		Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
646		parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
647		results?
648		Answer: [TODO]
649		Justification: [TODO]
650		Guidelines:
651		• The answer NA means that the paper does not include experiments.
652		• The experimental setting should be presented in the core of the paper to a level of detail
653		that is necessary to appreciate the results and make sense of them.
654		• The full details can be provided either with the code, in appendix, or as supplemental
655		material.
656	7.	Experiment statistical significance
657 658		Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?
659		Answer: [TODO]
660		Justification: [TODO]
661		Guidelines:
662		• The answer NA means that the paper does not include experiments.
663		• The authors should answer "Yes" if the results are accompanied by error bars, confi-
664		dence intervals, or statistical significance tests, at least for the experiments that support
665		the main claims of the paper.
666		• The factors of variability that the error bars are capturing should be clearly stated (for average train/test enlit initialization render drawing of some peremeter or everally
667 668		run with given experimental conditions).
669 670		• The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
671		• The assumptions made should be given (e.g., Normally distributed errors).
672		• It should be clear whether the error bar is the standard deviation or the standard error
673		of the mean.
674		• It is OK to report 1-sigma error bars, but one should state it. The authors should
675		preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
676		of Normality of errors is not verified.
677		• For asymmetric distributions, the authors should be careful not to show in tables or
679		error rates).
680		• If error bars are reported in tables or plots. The authors should explain in the text how
681		they were calculated and reference the corresponding figures or tables in the text.
682	8.	Experiments compute resources
683		Question: For each experiment, does the paper provide sufficient information on the com-
684		puter resources (type of compute workers, memory, time of execution) needed to reproduce
685		the experiments?
686		Answer: [TODO]
687		Justification: [TODO]
688		Guidelines:
689		• The answer NA means that the paper does not include experiments.
690		• The paper should indicate the type of compute workers CPU or GPU, internal cluster,
691		or cloud provider, including relevant memory and storage.

692 693	• The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
694	• The paper should disclose whether the full research project required more compute
695	than the experiments reported in the paper (e.g., preliminary or failed experiments that
696	didn't make it into the paper).
697	9. Code of ethics
698 699	Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
700	Answer: [TODO]
701	Justification: [TODO]
702	Guidelines:
703	• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
704	• If the authors answer No, they should explain the special circumstances that require a
705	deviation from the Code of Ethics.
706 707	• The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
708	10. Broader impacts
709 710	Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?
711	Answer: [TODO]
712	Justification: [TODO]
713	Guidelines:
714	• The answer NA means that there is no societal impact of the work performed.
715	• If the authors answer NA or No, they should explain why their work has no societal
716	impact or why the paper does not address societal impact.
717	• Examples of negative societal impacts include potential malicious or unintended uses
718	(e.g., disinformation, generating fake profiles, surveillance), fairness considerations
719	(e.g., deployment of technologies that could make decisions that unfairly impact specific
720	<ul> <li>The conference expects that many papers will be foundational research and not field</li> </ul>
722	to particular applications, let alone deployments. However, if there is a direct path to
723	any negative applications, the authors should point it out. For example, it is legitimate
724	to point out that an improvement in the quality of generative models could be used to
725	generate deepfakes for disinformation. On the other hand, it is not needed to point out
726	that a generic algorithm for optimizing neural networks could enable people to train
727	The outhors should consider possible horms that could arise when the technology is
728	• The authors should consider possible names that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the
730	technology is being used as intended but gives incorrect results, and harms following
731	from (intentional or unintentional) misuse of the technology.
732	• If there are negative societal impacts, the authors could also discuss possible mitigation
733	strategies (e.g., gated release of models, providing defenses in addition to attacks,
734	mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
735	feedback over time, improving the efficiency and accessibility of ML).
736	11. Safeguards
737	Question: Does the paper describe safeguards that have been put in place for responsible
738	release of data or models that have a high risk for misuse (e.g., pretrained language models,
739	mage generators, or scraped datasets):
740	
741	
742	Guidelines:
743	• The answer NA means that the paper poses no such risks.

744 745 746 747		• Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
748 749		• Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
750 751 752		• We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.
753	12.	Licenses for existing assets
754 755 756		Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?
757		Answer: [TODO]
758		Justification: [TODO]
759		Guidelines:
760		• The answer NA means that the paper does not use existing assets.
761		• The authors should cite the original paper that produced the code package or dataset.
762		• The authors should state which version of the asset is used and, if possible, include a
763		URL.
764		• The name of the license (e.g., CC-BY 4.0) should be included for each asset.
765		• For scraped data from a particular source (e.g., website), the copyright and terms of
766		service of that source should be provided.
767		• If assets are released, the license, copyright information, and terms of use in the
768		package should be provided. For popular datasets, paperswithcode.com/datasets
769		license of a dataset
770		• For existing datasets that are re-packaged both the original license and the license of
772		the derived asset (if it has changed) should be provided.
773 774		• If this information is not available online, the authors are encouraged to reach out to the asset's creators.
775	13.	New assets
776 777		Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?
778		Answer: [TODO]
779		Justification: [TODO]
780		Guidelines:
781		• The answer NA means that the paper does not release new assets.
782		• Researchers should communicate the details of the dataset/code/model as part of their
783		submissions via structured templates. This includes details about training, license,
784		limitations, etc.
785		• The paper should discuss whether and how consent was obtained from people whose
786		asset is used.
787 788		• At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
789	14.	Crowdsourcing and research with human subjects
790		Question: For crowdsourcing experiments and research with human subjects, does the paper
791 792		include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?
702		Answer: [TODO]
704		Institution: [TODO]
/ 94		

795	Guidelines:
796 797	• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
798 799 800 801 802	<ul> <li>Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.</li> <li>According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data</li> </ul>
803	collector.
804 15 805	. Institutional review board (IRB) approvals or equivalent for research with human subjects
806 807 808 809	Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?
810	Answer: [TODO]
811	Justification: [TODO]
812	Guidelines:
813 814	• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
815 816 817	• Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
818 819 820	• We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
821 822	• For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.
823 16	Declaration of LLM usage
824 825 826 827	Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.
828	Answer: [TODO]
829	Justification: [TODO]
830	Guidelines:
831 832	• The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
833 834	• Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.