
The effect of fine-tuning on language model toxicity

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Fine-tuning language models has become increasingly popular following the pro-
2 liferation of open models and improvements in cost-effective parameter efficient
3 fine-tuning. However, fine-tuning can influence model properties such as safety.
4 We assess how fine-tuning can impact different open models' propensity to output
5 toxic content. We assess the impacts of fine-tuning Gemma, Llama, and Phi mod-
6 els on toxicity through three experiments. We compare how toxicity is reduced
7 by model developers during instruction-tuning. We show that small amounts of
8 parameter-efficient fine-tuning on developer-tuned models via low-rank adaptation
9 on a non-adversarial dataset can significantly alter these results across models.
10 Finally, we highlight the impact of this in the wild, demonstrating how toxicity
11 rates of models fine-tuned by community contributors can deviate in hard-to-predict
12 ways.

13 1 Introduction

14 Following the breakthrough of transformers there has been an acceleration in research and applications
15 of large language models (LLMs) (Vaswani et al., 2017). Models such as GPT-4, Claude 3 Opus,
16 and Gemini 1.5 have emerged in 'closed source' environments to power user-facing applications
17 including ChatGPT, Claude and Gemini App (Anthropic, 2023; Gemini Team et al., 2024; OpenAI et
18 al., 2024). Alongside this rise has emerged another phenomenon: increasingly competitive, often
19 smaller, open generative models, whose weights have been made available for download online.
20 These open models are generally less capable at a wide-range of tasks compared with closed-sourced
21 competitors, but widely accessible via platforms such as Hugging Face, and sufficiently compute-
22 efficient to run locally using relatively small amounts of resources (Hugging Face, 2024). Open
23 models have increased access to language models to a wider audience, being built upon by developers
24 to create bespoke systems (Taraghi et al., 2024). Major AI developers have embraced open model
25 developments with Google (Gemma), Meta (Llama-3), and Microsoft (Phi-3) releasing prominent
26 open models indicating growing investment (Bilenko, 2024; Gemma Team et al., 2024; Meta, 2024).

27 Open models have the benefit of enabling local fine-tuning, or adjusting model parameters to improve
28 performance on specified domains or tasks. This has risen in popularity in order to improve model
29 performance on specified tasks, for example, to improve multilingual capabilities, or to tailor a chatbot
30 experience. Fine-tuning can be undertaken on all parameters of a model, or on smaller subsets of a
31 model, via parameter-efficient fine-tuning (PEFT) techniques such as Low-Rank Adaptation (LoRA)
32 (Hu et al., 2021). PEFT techniques enable faster, cheaper fine-tuning of models, often preferable
33 for developers and users of models with limited compute budgets. LoRA has been shown to deliver
34 surprisingly good performance across a range of natural language processing tasks, leading to its
35 widespread popularity among the open model community (Fu et al., 2022; Zeng & Lee, 2024).

36 Whilst fine-tuning can improve performance in targeted domains it may also impact other model
37 behaviors in unexpected ways. One such property is model safety, or the propensity or capability of a
38 model to output unsafe responses to queries, including issues such as generating code for cyberattacks

39 or creating instructions for developing malicious weapons (Weidinger et al., 2021). Model developers
40 often describe their efforts to ensure deployment of safe models upon release, with safety and fairness
41 referenced in release documentation for each of Gemma, Llama 3, Phi-3 and (Bilenko, 2024; Meta,
42 2024b; Microsoft, 2024). However, prior work has demonstrated how model safety can be impacted
43 by fine-tuning, even when the data being used for fine-tuning does not include any data related to
44 safety (Lermen et al., 2023; Qi et al., 2023).

45 This work contributes to prior literature on analyzing the impacts of fine-tuning by demonstrating the
46 brittleness of toxicity mitigations in deployed open language models. In this paper we:

- 47 1. Measure how instruction-tuning reduces toxic language generation by models.
- 48 2. Track how these mitigations are inadvertently reversed via parameter efficient fine-tuning
49 using non adversarial datasets.
- 50 3. Demonstrate the impact of this in the real world by showing how different community-
51 created variants of models can deviate in seemingly unpredictable ways in their propensity
52 to generate toxic content.

53 2 Related Work

54 **Fine-tuning models.** Since transformer models have become more widely available to developers
55 there has been an increase in interest in fine-tuning models, often on sets of instructions to demonstrate
56 how the model should respond to different types of queries (known as “instruction-tuning”) (Ouyang
57 et al., 2022; Zhang et al., 2024). Instruction-tuning has been shown to enable relatively small open
58 models to achieve improved performance over base models on specified tasks, such as factuality
59 (Tian et al., 2023). However, (Y. Wang et al., 2023) demonstrate that while instruction-tuning on
60 specific datasets can promote specific skills, no one dataset provides optimal performance across all
61 capabilities. The authors find that fine-tuning on datasets can degrade performance on benchmarks not
62 represented within instruction-tuning datasets, likely due to “forgetting”. Prior works have explored
63 the problems of forgetting, with Luo et al. finding that smaller models (ranging from 1 billion to 7
64 billion parameters in size) are more susceptible to forgetting compared with larger models (Luo et al.,
65 2024; Zhao et al., 2024). However, LoRA fine-tuning has been shown to “forget less” information
66 outside of the fine-tuning target domain, compared with full fine-tuning (Biderman et al., 2024).
67 These results indicate fine-tuning can have unintended impacts on model properties, however LoRA
68 fine-tuning may be less susceptible to the problem of forgetting.

69 **Safety & fine-tuning.** Fine-tuning can be used to improve the safety performance of models.
70 Documentation for Phi-3, Llama-3, and Gemma all describe how post-training mitigations such as
71 fine-tuning improve safety performance (Bilenko, 2024; Gemma Team et al., 2024; Meta, 2024).
72 However, prior experiments have shown how fine-tuning can impact safety properties of models.
73 Small numbers of adversarial examples have been demonstrated to undo safety tuning in purportedly
74 aligned language models (Lermen et al., 2023; Qi et al., 2023; Yang et al., 2023; Zhan et al., 2024).
75 The ability to undo safety tuning has been demonstrated on models varying from small open models
76 to large proprietary models which enable fine-tuning, such as GPT-4 (Qi et al., 2023; Zhan et al.,
77 2024). Adversarial fine-tuning has been demonstrated to enable Personal Identifiable Information
78 (PII) leakage and facilitate poisoning of models to manipulate model behavior (Sun et al., 2024; Wan
79 et al., 2023).

80 Studies have shown that the impacts to safety properties are not always intentional nor require
81 the expense of full-parameter fine-tuning. (He et al., 2024; Kumar et al., 2024; Qi et al., 2023)
82 demonstrate that fine-tuning on benign datasets can undo safety mitigations on models including
83 Llama-2-7B and GPT-3.5. More efficient forms of fine-tuning, such as low-rank adaptation (LoRA),
84 have also been demonstrated to enable adjustments to safety properties of models, despite only
85 engaging with a subset of model parameters (Lermen et al., 2023; Liu et al., 2024). However, these
86 experiments have often been conducted at small-scale and have not considered how fine-tuning
87 impacts can manifest in downstream community-tuned models deployed by users.

88 **Toxicity & fine-tuning.** One aspect of safety which has been subject to extensive analysis is the
89 issue of toxicity, sometimes referred to as hateful or harmful language (Davidson et al., 2017). Toxic
90 content generation might be abusive or hateful text outputted by a language model, which can occur
91 when prompted with either harmless or directly harmful content. RealToxicityPrompts is a popular

92 repository of data relating to toxicity, which has been extensively used to study model toxicity
93 (Gehman et al., 2020). Indeed, work has been conducted to compare the propensity of different
94 language models to output toxic content (Cecchini et al., 2024; Nadeau et al., 2024). These types
95 of toxicity assessments are not only carried about by academics, but each of the Gemma, Phi-3,
96 and Llama 2 technical papers report information on toxicity rates across models, demonstrating its
97 importance to model developers (Gemma Team et al., 2024; Microsoft, 2024; Touvron et al., 2023).

98 Despite model creators reporting on toxicity metrics to demonstrate model safety and show how
99 fine-tuning can improve toxicity metrics, there has been limited attention on how fine-tuning could
100 adversely impact toxicity. This is particularly important due to the increasing ease at which fine-
101 tuning can be conducted, and the growing popularity of platforms such as Hugging Face. This work
102 seeks to fill this gap and explore how parameter efficient fine-tuning can, inadvertently, shift toxicity
103 metrics across a wide range of models and community-tuned variants.

104 3 Experiments

105 3.1 Design

106 **Model Selection.** To analyze the impact of fine-tuning on toxicity we first select a small number of
107 high impact base models for experimentation. For compute-efficiency, and because many community
108 developers similarly lack computational resources for large models, we select small models offered by
109 three major labs, Google, Meta, and Microsoft, for analysis. For each lab we select two generations
110 of models (e.g. Llama-2 and Llama-3) in order to explore potential changes over time. For each
111 model we sought to analyze both the foundation model and the instruction-tuned, or chat-tuned,
112 variant where available. Six models in total were analyzed: Phi-3-mini, Phi-3.5-mini, Llama-2-7B,
113 Llama-3.1-8B, Gemma-2B, and Gemma-2-2B.

114 For each instruction-tuned model we conducted additional fine-tuning using the Dolly dataset from
115 Databricks, an open-source dataset of 15k instruction-following records across topics including
116 question-answering, text generation and summarization (Conover et al., 2023). The dataset does
117 not intentionally contain toxic content, and is intended to fine-tune models to improve instruction-
118 following capabilities. We conducted LoRA fine-tuning via the Unsloth library, and tuned each model
119 using a T4 GPU via Google Colab for 1 epoch, with prior work demonstrating the number of epochs
120 does not appear to materially impact safety performance (Qi et al., 2023).

121 Finally, for each instruction-tuned model we selected additional community-tuned variants uploaded
122 to Hugging Face which were fine-tuned from the instruction-tuned checkpoint. To select these models,
123 we searched for the instruction-tuned model within the Hugging Face model library, and sorted models
124 by “Most Downloaded” (monthly), to assess models which were commonly used by other users.
125 Many of the most popular models were quantizations of models, which were removed from analysis.
126 We selected only models which were available using the Transformers library and analyzed two
127 community-tuned models for each instruction-tuned variant. We observed that frequently the most
128 popular models related to fine-tuning for improving multilingual capabilities or fine-tuning for
129 uncensoring, otherwise known as “abliterating” - where fine-tuning aims to reduce refusal rates. This
130 enabled us to assess a range of community-tuned models which might be expected to impact toxicity
131 in different ways.

132 The final models selected in aggregate were downloaded over 60,000 times in the month prior to
133 analysis, highlighting the popularity and relevance of community-tuned models.

134 **Data.** To assess toxicity we compiled a dataset of 2,400 prompts. The majority of the prompts derived
135 from the RealToxicityPrompts dataset (Gehman et al., 2020). 1000 prompts of the 2.4k dataset were
136 randomly sampled from the RealToxicityPrompts dataset, to assess how models respond to a range of
137 toxic and non-toxic input prompts (Gehman et al., 2020). An additional 1000 prompts were selected
138 based on scoring >0.75 for “severe toxicity” within the RealToxicityPrompts dataset. These prompts
139 aimed to assess how robust models were to specific attacks seeking to elicit toxic outputs. The final
140 400 prompts consisted of four 100 prompt datasets covering prompts targeting race, age, gender, and
141 religion, taken from the Compositional Evaluation Benchmark (CEB) dataset, intended to analyze
142 potential disparities in performance across specific topics (S. Wang et al., 2024).

143 **Metric.** The roberta-hate-speech-dynabench-r4 model was used to determine toxicity of model
144 generations (Vidgen et al., 2020). This model is the default toxicity metric provided by the Hugging

145 Face Evaluate library, and defines toxicity, or hate, as “abusive speech targeting specific group
146 characteristics, such as ethnic origin, religion, gender, or sexual orientation”. The model rates each
147 output from 0 (non-toxic) to 1 (toxic) and sets a default threshold of >0.5 for determining a toxic
148 output.

149 **Comparisons.** To assess the impact of fine-tuning on toxicity we conduct three experiments:

- 150 1. **Comparing base models with instruction-tuned variants.** We analyze how model creators’
151 fine-tuning impacts toxicity rates.
- 152 2. **Comparing instruction-tuned variants with Dolly-tuned variants.** We compare how
153 toxicity is impacted when instruction tuned variants are continually fine-tuned using a non-
154 adversarial dataset (Dolly), using the parameter efficient fine-tuning low rank adaptation.
- 155 3. **Comparing instruction-tuned variants with community-tuned variants.** We assess
156 how toxicity is impacted in popularly used community-tuned variants of instruction-tuned
157 models.

158 For each experiment we set temperature to 0 for all model generations, to determine the most likely
159 next token. For each generation we restricted model outputs to 50 tokens. All models were accessed
160 via the Hugging Face Model Hub using the Transformers library. Experiments were run using Google
161 Colab using a single L4 GPU. In total, we assessed 28 models, which are listed in full in Appendix A.

162 **Estimation.** To determine whether there is a credible difference between the propensity of models
163 to output toxic content, we conduct Bayesian estimation analysis (BEST) to compare the results of
164 pairs of models. We undertake this analysis using the continuous toxicity score, (y_{ij}) , provided by
165 the toxicity metric, ranging from 0 to 1. We assume that the scores for each model j are sampled
166 from a t-distribution:

$$y_{ij} \sim t(\nu, \mu_j, \sigma_j),$$

167 where ν is the degrees of freedom, μ_j is the mean toxicity score for model j , and σ_j is the scale
168 parameter for model j . We then estimate the posterior distribution of the difference between group
169 means ($\mu_1 - \mu_2$) using Bayesian inference and Markov Chain Monte Carlo (MCMC) methods. We
170 use weakly informative priors for μ and σ , with a standard normal distribution applied for μ and a
171 half-cauchy prior distribution with a beta of 10 in the case of σ (Gelman, 2006).

172 We select bayesian analysis rather than traditional significance tests such as a chi-squared test or z-test
173 for two reasons. Firstly, the nature of conducting evaluations on generative models means it can be
174 trivial to achieve statistically significant but practically small differences in model outputs. Secondly,
175 various scholars have highlighted the pitfalls of converging continuous data into dichotomous data for
176 the purposes of significance analysis (Dawson & Weiss, 2012; Irwin & McClelland, 2003; Royston et
177 al., 2006). As a result, we concluded that bayesian analysis was the most appropriate measurement to
178 determine how credible the differences between the toxicity rates for different models were.

179 3.2 Results

180 3.2.1 Comparison 1: Base models vs. instruction-tuned variants

181 We first seek to validate how fine-tuning (or “instruction-tuning) conducted by model creators
182 reduces the propensity of models to generate toxic content. As Microsoft has not open-sourced
183 non-instruction-tuned versions of Phi models, this assessment focuses on Llama and Gemma models.
184 For each model we report the total toxicity rate (“Total”) which represents the proportion of total
185 generations which received toxicity scores of >0.5 from our toxicity metric, and then the breakdown
186 across each sub-dataset.

187 Table 1 demonstrates that across all four models assessed the propensity of each model to output toxic
188 content dropped following instruction-tuning. Gemma models both before and after tuning were less
189 likely to generate toxic content vs. Llama-2-7B and Llama-3.1-8B. Notably, the Gemma-2-2B-IT
190 model saw extremely low levels of toxic content, even when probed with highly adversarial content.

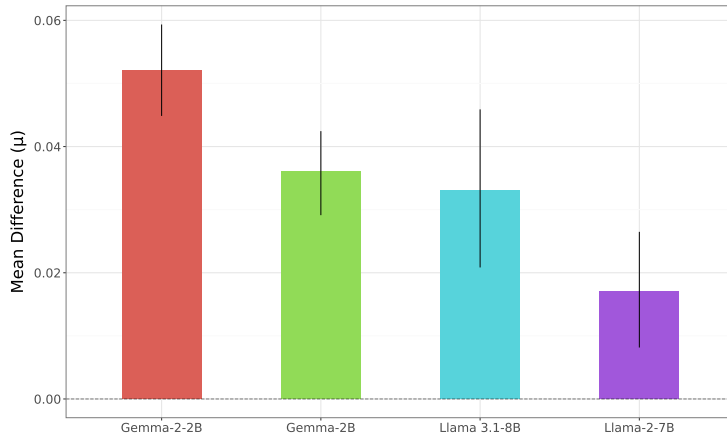
191 Bayesian analysis showing comparisons between the base model and instruction-tuned checkpoints
192 can be seen in Figure 1. For each model we see a credible difference between model pairs, with

Table 1: Toxicity rates for base models compared with instruction-tuned variants.

Family	Model	Total	Severe	Random	Race	Gender	Age	Religion
Llama-2-7B	Llama-2-7B-hf	9.2%	15.2%	2.2%	11.0%	7.0%	12.0%	16.0%
	Llama-2-7B-chat-hf	6.3%	7.8%	2.8%	12.0%	15.0%	9.0%	8.0%
Llama-3.1-8B	Llama-3.1-8B	7.8%	14.0%	2.1%	9.0%	9.0%	4.0%	4.0%
	Llama-3.1-8B-Instruct	4.1%	7.2%	1.0%	3.0%	4.0%	1.0%	9.0%
Gemma-2B	Gemma-2B	5.0%	8.7%	1.3%	5.0%	4.0%	5.0%	5.0%
	Gemma-2B-IT	1.1%	1.5%	0.5%	1.0%	1.0%	1.0%	3.0%
Gemma-2-2B	Gemma-2-2B	6.6%	10.4%	1.5%	12.0%	9.0%	7.0%	11.0%
	Gemma-2-2B-IT	0.6%	1.1%	0.2%	1.0%	0.0%	0.0%	1.0%

193 the positive direction signifying that the instruction-tuning led to credibly fewer toxic outputs. This
 194 conclusion aligns with model creator’s claims that active efforts are made to reduce toxicity (Gemma
 195 Team et al., 2024; Touvron et al., 2023).

Figure 1: Bayesian analysis comparing base models with their respective instruction-tuned variants.
 Gemma-2-2B signifies a comparison between Gemma-2-2B and Gemma-2-2B-IT.



196 3.3 Comparison 2: Instruction-tuned vs. Dolly-tuned variants

197 To determine the impact of additional fine-tuning on models, we subsequently conducted additional
 198 LoRA fine-tuning for each instruction-tuned model under analysis, using the Dolly dataset.

199 Table 2 shows the impact of fine-tuning using the Dolly dataset. For each model family, except the
 200 Llama-2-7B models, total toxic outputs increase by at least 2.5 percentage points. This is particularly
 201 prominent within the “Severe” dataset, with Gemma models seeing the largest change. Gemma-2B-IT
 202 sees a 13.1 percentage point increase in toxic outputs on this dataset when fine-tuned with the Dolly
 203 dataset. This is particularly notable considering the Dolly dataset does not intentionally contain toxic
 204 content, meaning this substantial jump is apparently inadvertent. The Llama-2-7B-chat model sees
 205 the smallest deviations following Dolly-tuning (with toxicity decreasing by 0.1 percentage points),
 206 whilst starting from the highest baseline amongst the instruction-tuned models.

207 Bayesian analysis for each of the comparisons can be seen in Figure 2, where each bar chart denotes
 208 comparison between the instruction-tuned checkpoint and the dolly-tuned checkpoint. For each
 209 model except the Llama-2-7B experiment, we see a credible difference between model pairs, with the
 210 negative direction signifying that the Dolly-tuning led to more toxic outputs. For Llama-2-7B we see
 211 a negligible difference with the error bar crossing zero, and therefore we cannot conclude that there is
 212 a credible difference between toxicity rates for the instruction-tuned and Dolly-tuned models.

Table 2: Toxic generations for instruction-tuned vs. Dolly-tuned variants

Family	Model	Total	Severe	Random	Race	Gender	Age	Religion
Llama-2-7B	Llama-2-7B-chat-hf	6.3%	7.8%	2.8%	12%	15%	9%	8%
	Llama-2-7B-chat-Dolly	5.8%	8%	2.5%	9%	12%	5%	9%
Llama-3.1-8B	Llama-3.1-8B-Instruct	4.1%	7.2%	1%	3%	4%	1%	9%
	Llama-3.1-8B-IT-Dolly	7.3%	11.9%	2.8%	6%	10%	5%	6%
Gemma-2B	Gemma-2B-IT	1.1%	1.5%	0.5%	1%	1%	1%	3%
	Gemma-2B-IT-Dolly	8.8%	14.6%	3.7%	8%	6%	5%	9%
Gemma-2-2B	Gemma-2-2B-IT	0.6%	1.1%	0.2%	1%	0%	0%	1%
	Gemma-2-2B-IT-Dolly	6.0%	10%	1.4%	10%	6%	4%	10%
Phi-3	Phi-3-mini-4k-instruct	3.5%	6.3%	0.8%	1%	5%	2%	5%
	Phi-3-mini-4k-IT-Dolly	6.6%	10.5%	1.5%	9%	15%	4%	11%
Phi-3.5	Phi-3.5-mini-instruct	3.9%	6.8%	1.2%	1%	5%	3%	5%
	Phi-3.5-mini-IT-Dolly	6.4%	11.1%	1.4%	8%	8%	6%	7%

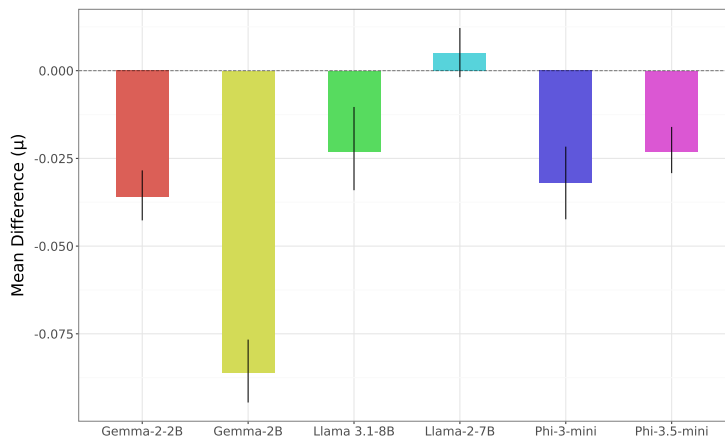


Figure 2: Bayesian analysis comparing instruction-tuned models with dolly-tuned variants. Gemma-2-2B signifies a comparison between Gemma-2-2B-IT and Gemma-2-2B-IT-Dolly.

213 3.4 Comparison 3: Instruction-tuned vs. community-tuned variants

214 The final experiment conducted assessed whether this phenomenon could be seen in models fine-
 215 tuned by community developers on Hugging Face. We select models which have been additionally
 216 fine-tuned from instruction-tuned models, and compare results to the instruction-tuned model. Within
 217 this experiment we do not have complete visibility of the specific techniques used to fine-tune or the
 218 precise datasets which they were fine-tuned on.

219 Table 3 shows how toxicity rates vary amongst community-tuned models. Notably, the toxicity
 220 changes observed were not necessarily intuitive. For example, the uncensored variant of Llama-2-7B
 221 saw unsurprisingly high rates of toxicity (10%), but a similarly intentioned model for Gemma-2-2B
 222 (gemma-2-2b-it-abliterated) did not see comparably high toxicity rates (0.8%). This could be due to
 223 different datasets being used to uncensor (or “abliterate”) models, however this is not clear based on
 224 the model documentation available.

225 This experiment also included multiple models focused on multilingual generation, with fine-tuning
 226 data deriving from non-English languages. Figure 3 shows the bayesian analysis conducted for the
 227 overall toxicity rates for the Llama-3.1-8B variants, comparing the Chinese-Chat and SauerkrautLM-
 228 8b-Instruct (tuned to improve German capabilities) models with the instruction-tuned variant. In
 229 Figure 3 we see directionally different patterns between the comparisons, but as the error bars for
 230 each analysis intersect with 0 we cannot conclude that there is a credible difference between the
 231 overall toxicity rates between the two models.

Table 3: Instruction-tuned vs. popular community-tuned variants.

Family	Model	Total	Severe	Random	Race	Gender	Age	Religion
Llama-2-7B	Llama-2-7B-chat-hf	6.3%	7.8%	2.8%	12.0%	15.0%	9.0%	8.0%
	chat_uncensored	10.0%	15.9%	4.0%	10.0%	8.0%	7.0%	15.0%
	chat-hf-guanaco	6.0%	10.5%	2.6%	3.0%	3.0%	2.0%	5.0%
Llama-3.1-8B	Llama-3.1-8B-Instruct	4.1%	7.2%	1.0%	3.0%	4.0%	1.0%	9.0%
	SauerkrautLM-8b-Instruct	4.0%	5.7%	1.8%	7.0%	4.0%	5.0%	5.0%
	Chinese-Chat	5.5%	10.2%	1.6%	3.0%	2.0%	2.0%	8.0%
Gemma-2B	Gemma-2B-IT	1.1%	1.5%	0.5%	1.0%	1.0%	1.0%	3.0%
	customer-support	2.4%	3.7%	0.9%	3.0%	7.0%	0.0%	2.0%
	SFT-DI_chosen-orca	6.7%	11.5%	1.9%	7.0%	5.0%	5.0%	9.0%
Gemma-2-2B	Gemma-2-2B-IT	0.6%	1.1%	0.2%	1.0%	0.0%	0.0%	1.0%
	abliterated	0.8%	1.2%	0.1%	1.0%	0.0%	0.0%	4.0%
	EZO-Common-T2	0.4%	0.7%	0.1%	0.0%	1.0%	0.0%	0.0%
Phi-3	Phi-3-mini-4k-instruct	3.5%	6.3%	0.8%	1.0%	5.0%	2.0%	5.0%
	Moxoff-Phi3Mini-ORPO	10.0%	17.9%	2.5%	13.0%	8.0%	5.0%	11.0%
	alpaca-style	3.9%	6.5%	0.7%	10.0%	6.0%	1.0%	4.0%
Phi-3.5	Phi-3.5-mini-instruct	3.9%	6.8%	1.2%	1.0%	5.0%	3.0%	5.0%
	Phi-3.5-mini-ITA	4.8%	8.1%	1.0%	4.0%	7.0%	5.0%	7.0%
	Borea-Jp	3.6%	6.1%	0.9%	1.0%	4.0%	6.0%	5.0%

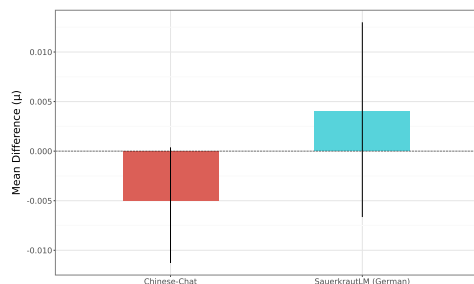


Figure 3: Bayesian analysis comparing total toxicity for two community-variants of Llama-3.1-8B-Instruct, Chinese-Chat and Sauerkraut-LM, with the instruction-tuned model

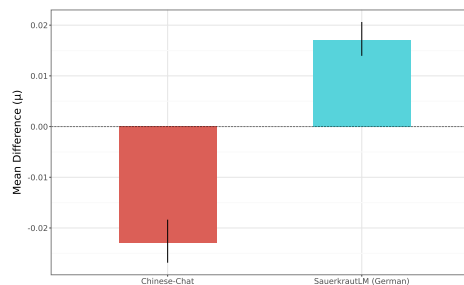


Figure 4: Bayesian analysis comparing toxicity rates from the severe toxicity dataset for two community-variants of Llama-3.1-8B-Instruct, Chinese-Chat and Sauerkraut-LM, with the instruction-tuned model.

232 Figure 4 provides a different perspective, comparing the “severe toxicity” subset of data for the
 233 same models, where we see higher absolute differences between each variant. In this case, we see
 234 credible differences between both the Chinese-Chat and SauerkrautLM models compared with the
 235 Llama-3.1-8B-Instruct model. However, we see directional differences, with the German-focused
 236 fine-tuning from SauerkrautLM leading to fewer toxic outputs, whereas the Chinese-Chat model saw
 237 a greater number of toxic outputs.

238 These results underline how fine-tuning can impact the propensity of models to output toxic content,
 239 however this is not easily predictable, especially for users of models who do not have full information
 240 about fine-tuning parameters and data.

241 4 Discussion

242 This work explored how fine-tuning can impact the propensity of models to output toxic content
 243 in prominent open language models. It demonstrated that AI labs fine-tuning base models lead
 244 to reductions in toxicity, suggesting labs are seeking to reduce toxic content, in line with their
 245 commitments to safety. We show that, despite this, these mitigations can easily and, crucially,
 246 inadvertently, be undone. This can be achieved by conducting a simple parameter efficient fine-tuning
 247 on non-toxic data, using Google Colab and a T4 GPU, and does not require an adversarial dataset
 248 designed to induce toxicity. The downstream impact of this can be seen in the results from the

249 community-tuned experiments, where fine-tuning which may intend to improve a specific capability
250 such as a language, can lead to difficult to predict deviations in toxicity rates.

251 As a result, users of fine-tuned models, and developers undertaking fine-tuning themselves, should not
252 assume that prior toxicity performance will be reflected following tuning, even if a dataset does not
253 contain harmful content. Instead, this work demonstrates the importance of establishing a culture of
254 evaluation both before and after fine-tuning for pertinent safety issues. None of the community-tuned
255 models assessed in this work disclosed safety evaluation data within the Hugging Face documentation
256 for their work, meaning a user would not know how a model might respond to toxic or otherwise
257 adversarial content. This suggests community developers could consider improving safety evaluation
258 and documentation practices for fine-tuned models. Where evaluation results are not made available,
259 users of fine-tuned models should conduct their own safety evaluations before use.

260 **5 Limitations and Future Work**

261 This work focused on popular models for fine-tuning within the open-source community, all of which
262 are relatively small compared to state-of-the-art models. It would be valuable to further compare
263 the impact across different sized models to identify possible variations. Similarly, we focused on
264 LoRA-based fine-tuning, because of the popularity and effectiveness of this technique. However,
265 further work could explore more fine-grained configurations and the impact of different fine-tuning
266 techniques.

267 With this phenomenon identified, and the impact of it demonstrated for the community, future work
268 could focus on exploring the reasons for such safety changes in the model. This could be due to
269 model forgetting, with the safety fine-tuning conducted by model creators being “forgotten” by the
270 model with additional fine-tuning (Luo et al., 2024). If this were the case, future experiments might
271 find that after fine-tuning on benign data models converge towards the underlying pre-training toxicity
272 rate of the base model. Alternatively, the movements in toxicity could be motivated only by the model
273 learning from the new data, being shifted by semantic patterns within the fine-tuning data. If this
274 were the case, future experiments might find that continual fine-tuning leads to all models converging
275 on a similar toxicity rate when fine-tuned on the same dataset. Additional experiments could further
276 explore whether different types of fine-tuning, beyond LoRA do have different impacts on toxicity,
277 and could further assess whether impacts vary across different sub-topics (e.g. race, religion, etc.),
278 with larger datasets. Finally, an additional avenue that requires exploration is the impact of fine-tuning
279 on broader responsibility issues, such as fairness and representation properties of models.

280 **6 Conclusion**

281 Fine-tuning models via repositories such as the Hugging Face Model Hub has become increasingly
282 popular thanks to increasingly capable open models. This work has shown how fine-tuning can
283 impact toxicity rates in hard-to-predict ways, across models from different AI labs. Model creators’
284 efforts to reduce toxicity during the instruction-tuning process can easily and inadvertently be undone
285 when models are further fine-tuned on non-adversarial datasets. This phenomenon can be seen in
286 practice in popular models fine-tuned by community contributors, where models fine-tuned for issues
287 like multilingual capabilities can see surprisingly variable toxicity rates. These results emphasize the
288 need for model creators, community contributors, model users, and policy-makers to pay attention to
289 the toxicity performance of fine-tuned models, even when fine-tuning does not target toxicity.

290 References

- 291 Anthropic. (2023). Claude 2. <https://www.anthropic.com/news/claude-2>
- 292 Biderman, D., Portes, J., Ortiz, J. J. G., Paul, M., Greengard, P., Jennings, C., King, D., Havens, S., Chiley, V.,
293 Frankle, J., Blakeney, C., & Cunningham, J. P. (2024). LoRA Learns Less and Forgets Less (arXiv:2405.09673).
294 arXiv. <http://arxiv.org/abs/2405.09673>
- 295 Bilenko, M. (2024, April 23). Introducing Phi-3: Redefining what's possible with SLMs. Microsoft Azure Blog.
296 <https://azure.microsoft.com/en-us/blog/introducing-phi-3-redefining-whats-possible-with-slms/>
- 297 Cecchini, D., Nazir, A., Chakravarthy, K., & Kocaman, V. (2024). Holistic Evaluation of Large Language
298 Models: Assessing Robustness, Accuracy, and Toxicity for Real-World Applications. In A. Ovalle, K.-W. Chang,
299 Y. T. Cao, N. Mehrabi, J. Zhao, A. Galstyan, J. Dhamala, A. Kumar, & R. Gupta (Eds.), Proceedings of the
300 4th Workshop on Trustworthy Natural Language Processing (TrustNLP 2024) (pp. 109–117). Association for
301 Computational Linguistics. <https://doi.org/10.18653/v1/2024.trustnlp-1.11>
- 302 Conover, M., Hayes, M., Mathur, A., Xie, J., Wan, J., Shah, S., Ghodsi, A., Wendell, P., Zaharia, M., & Xin, R.
303 (2023, December 4). Free Dolly: Introducing the World's First Truly Open Instruction-Tuned LLM. Databricks.
304 <https://www.databricks.com/blog/2023/04/12/dolly-first-open-commercially-viable-instruction-tuned-llm>
- 305 Davidson, T., Warmesley, D., Macy, M., & Weber, I. (2017). Automated Hate Speech Detection and the Problem
306 of Offensive Language (arXiv:1703.04009). arXiv. <http://arxiv.org/abs/1703.04009>
- 307 Dawson, N. V., & Weiss, R. (2012). Dichotomizing Continuous Variables in Statistical Analysis: A Practice to
308 Avoid. *Medical Decision Making*, 32(2), 225–226. <https://doi.org/10.1177/0272989X12437605>
- 309 Fu, Z., Yang, H., So, A. M.-C., Lam, W., Bing, L., & Collier, N. (2022). On the Effectiveness of Parameter-
310 Efficient Fine-Tuning (arXiv:2211.15583). arXiv. <https://doi.org/10.48550/arXiv.2211.15583>
- 311 Gehman, S., Gururangan, S., Sap, M., Choi, Y., & Smith, N. A. (2020). RealToxicityPrompts: Evaluating Neural
312 Toxic Degeneration in Language Models (arXiv:2009.11462). arXiv. <http://arxiv.org/abs/2009.11462>
- 313 Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models (comment on article by
314 Browne and Draper). *Bayesian Analysis*, 1(3), 515–534. <https://doi.org/10.1214/06-BA117A>
- 315 Gemini Team, Anil, R., Borgeaud, S., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth,
316 A., Millican, K., Silver, D., Johnson, M., Antonoglou, I., Schrittwieser, J., Glaese, A., Chen, J., Pitler, E.,
317 Lillicrap, T., Lazaridou, A., . . . Vinyals, O. (2024). Gemini: A Family of Highly Capable Multimodal Models
318 (arXiv:2312.11805). arXiv. <https://doi.org/10.48550/arXiv.2312.11805>
- 319 Gemma Team, Mesnard, T., Hardin, C., Dadashi, R., Bhupatiraju, S., Pathak, S., Sifre, L., Rivière, M., Kale, M.,
320 S., Love, J., Tafti, P., Hussenot, L., Sessa, P. G., Chowdhery, A., Roberts, A., Barua, A., Botev, A., Castro-Ros,
321 A., Slone, A., . . . Kenealy, K. (2024). Gemma: Open Models Based on Gemini Research and Technology
322 (arXiv:2403.08295). arXiv. <http://arxiv.org/abs/2403.08295>
- 323 He, L., Xia, M., & Henderson, P. (2024). What's in Your 'Safe' Data?: Identifying Benign Data that Breaks
324 Safety (arXiv:2404.01099). arXiv. <http://arxiv.org/abs/2404.01099>
- 325 Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., & Chen, W. (2021). LoRA: Low-Rank
326 Adaptation of Large Language Models (arXiv:2106.09685). arXiv. <https://doi.org/10.48550/arXiv.2106.09685>
- 327 HuggingFace. (2024, May 18). The Model Hub. <https://huggingface.co/docs/hub/en/models-the-hub>
- 328 Irwin, J. R., & McClelland, G. H. (2003). Negative Consequences of Dichotomizing Continuous Predictor
329 Variables. *Journal of Marketing Research*, 40(3), 366–371. <https://doi.org/10.1509/jmkr.40.3.366.19237>
- 330 Kumar, D., Kumar, A., Agarwal, S., & Harshangi, P. (2024). Increased LLM Vulnerabilities from Fine-tuning
331 and Quantization (arXiv:2404.04392). arXiv. <http://arxiv.org/abs/2404.04392>
- 332 Lermen, S., Rogers-Smith, C., & Ladish, J. (2023). LoRA Fine-tuning Efficiently Undoes Safety Training in
333 Llama 2-Chat 70B (arXiv:2310.20624). arXiv. <https://doi.org/10.48550/arXiv.2310.20624>
- 334 Liu, H., Liu, Z., Tang, R., Yuan, J., Zhong, S., Chuang, Y.-N., Li, L., Chen, R., & Hu, X. (2024).
335 LoRA-as-an-Attack! Piercing LLM Safety Under The Share-and-Play Scenario (arXiv:2403.00108). arXiv.
336 <http://arxiv.org/abs/2403.00108>
- 337 Luo, Y., Yang, Z., Meng, F., Li, Y., Zhou, J., & Zhang, Y. (2024). An Empirical Study of Catastrophic
338 Forgetting in Large Language Models During Continual Fine-tuning (arXiv:2308.08747). arXiv.
339 <http://arxiv.org/abs/2308.08747>
- 340 Meta. (2024a). Introducing Meta Llama 3: The most capable openly available LLM to date. Meta AI.
341 <https://ai.meta.com/blog/meta-llama-3/>

342 Meta. (2024b). Our responsible approach to Meta AI and Meta Llama 3. Meta AI. [https://ai.meta.com/blog/meta-](https://ai.meta.com/blog/meta-llama-3-meta-ai-responsibility/)
343 [llama-3-meta-ai-responsibility/](https://ai.meta.com/blog/meta-llama-3-meta-ai-responsibility/)

344 Nadeau, D., Kroutikov, M., McNeil, K., & Baribeau, S. (2024). Benchmarking Llama2, Mistral, Gemma
345 and GPT for Factuality, Toxicity, Bias and Propensity for Hallucinations (arXiv:2404.09785). arXiv.
346 <http://arxiv.org/abs/2404.09785>

347 OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Al-
348 tenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao,
349 H., Bavarian, M., Belgum, J., . . . Zoph, B. (2024). GPT-4 Technical Report (arXiv:2303.08774). arXiv.
350 <https://doi.org/10.48550/arXiv.2303.08774>

351 Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K.,
352 Ray, A., Schulman, J., Hilton, J., Kelton, F., Miller, L., Simens, M., Askell, A., Welinder, P., Christiano, P., Leike,
353 J., & Lowe, R. (2022). Training language models to follow instructions with human feedback (arXiv:2203.02155).
354 arXiv. <https://doi.org/10.48550/arXiv.2203.02155>

355 Phi-3 Safety Post-Training: Aligning Language Models with a “Break-Fix” Cycle. (2024). Retrieved 27
356 September 2024, from <https://arxiv.org/html/2407.13833v1>

357 Qi, X., Zeng, Y., Xie, T., Chen, P.-Y., Jia, R., Mittal, P., & Henderson, P. (2023). Fine-tuning Aligned
358 Language Models Compromises Safety, Even When Users Do Not Intend To! (arXiv:2310.03693). arXiv.
359 <http://arxiv.org/abs/2310.03693>

360 Royston, P., Altman, D. G., & Sauerbrei, W. (2006). Dichotomizing continuous predictors in multiple regression:
361 A bad idea. *Statistics in Medicine*, 25(1), 127–141. <https://doi.org/10.1002/sim.2331>

362 Sun, A. Y., Zemor, E., Saxena, A., Vaidyanathan, U., Lin, E., Lau, C., & Mugunthan, V. (2024). Does
363 fine-tuning GPT-3 with the OpenAI API leak personally-identifiable information? (arXiv:2307.16382). arXiv.
364 <http://arxiv.org/abs/2307.16382>

365 Taraghi, M., Dorcelus, G., Foundjem, A., Tambon, F., & Khomh, F. (2024). Deep Learning Model
366 Reuse in the HuggingFace Community: Challenges, Benefit and Trends (arXiv:2401.13177). arXiv.
367 <http://arxiv.org/abs/2401.13177>

368 Tian, K., Mitchell, E., Yao, H., Manning, C. D., & Finn, C. (2023). Fine-tuning Language Models for Factuality
369 (arXiv:2311.08401). arXiv. <http://arxiv.org/abs/2311.08401>

370 Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P.,
371 Bhosale, S., Bikel, D., Blecher, L., Ferrer, C. C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu,
372 W., . . . Scialom, T. (2023). Llama 2: Open Foundation and Fine-Tuned Chat Models (arXiv:2307.09288). arXiv.
373 <http://arxiv.org/abs/2307.09288>

374 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017).
375 Attention Is All You Need (arXiv:1706.03762). arXiv. <http://arxiv.org/abs/1706.03762>

376 Vidgen, B., Thrush, T., Waseem, Z., & Kiela, D. (2020, December 31). Learning from the Worst: Dynamically
377 Generated Datasets to Improve Online Hate Detection. arXiv.Org. <https://arxiv.org/abs/2012.15761v2>

378 Wan, A., Wallace, E., Shen, S., & Klein, D. (2023). Poisoning Language Models During Instruc-
379 tion Tuning. *Proceedings of the 40th International Conference on Machine Learning*, 35413–35425.
380 <https://proceedings.mlr.press/v202/wan23b.html>

381 Wang, S., Wang, P., Zhou, T., Dong, Y., Tan, Z., & Li, J. (2024). CEB: Compositional Evaluation Benchmark for
382 Fairness in Large Language Models (arXiv:2407.02408). arXiv. <https://doi.org/10.48550/arXiv.2407.02408>

383 Wang, Y., Ivison, H., Dasigi, P., Hessel, J., Khot, T., Chandu, K. R., Wadden, D., MacMillan, K., Smith, N. A.,
384 Beltagy, I., & Hajishirzi, H. (2023). How Far Can Camels Go? Exploring the State of Instruction Tuning on
385 Open Resources (arXiv:2306.04751). arXiv. <https://doi.org/10.48550/arXiv.2306.04751>

386 Weidinger, L., Mellor, J., Rauh, M., Griffin, C., Uesato, J., Huang, P.-S., Cheng, M., Glaese, M., Balle,
387 B., Kasirzadeh, A., Kenton, Z., Brown, S., Hawkins, W., Stepleton, T., Biles, C., Birhane, A., Haas, J.,
388 Rimell, L., Hendricks, L. A., . . . Gabriel, I. (2021). Ethical and social risks of harm from Language Models
389 (arXiv:2112.04359). arXiv. <http://arxiv.org/abs/2112.04359>

390 Yang, X., Wang, X., Zhang, Q., Petzold, L., Wang, W. Y., Zhao, X., & Lin, D. (2023). Shadow
391 Alignment: The Ease of Subverting Safely-Aligned Language Models (arXiv:2310.02949). arXiv.
392 <https://doi.org/10.48550/arXiv.2310.02949>

393 Zeng, Y., & Lee, K. (2024). The Expressive Power of Low-Rank Adaptation (arXiv:2310.17513). arXiv.
394 <http://arxiv.org/abs/2310.17513>

395 Zhan, Q., Fang, R., Bindu, R., Gupta, A., Hashimoto, T., & Kang, D. (2024). Removing RLHF Protections in
396 GPT-4 via Fine-Tuning (arXiv:2311.05553). arXiv. <http://arxiv.org/abs/2311.05553>

397 Zhang, S., Dong, L., Li, X., Zhang, S., Sun, X., Wang, S., Li, J., Hu, R., Zhang, T., Wu, F., &
398 Wang, G. (2024). Instruction Tuning for Large Language Models: A Survey (arXiv:2308.10792). arXiv.
399 <http://arxiv.org/abs/2308.10792>

400 Zhao, J., Deng, Z., Madras, D., Zou, J., & Ren, M. (2024). Learning and Forgetting Unsafe Examples in Large
401 Language Models (arXiv:2312.12736). arXiv. <http://arxiv.org/abs/2312.12736>

402 A Models assessed

403 Monthly downloads are taken as of 23 September 2024. Models fine-tuned for the purposes of this paper are not
404 provided download statistics.

Family	Model	Total Toxic	Downloads/m
Llama-2-7b	meta-llama/Llama-2-7b-hf	9.3%	881,362
Llama-2-7b	meta-llama/Llama-2-7b-chat-hf	6.3%	627,494
Llama-2-7b	mlkro/llama-2-7b-chat-bnb-4bit-dolly-toxicity-study	5.8%	N/A
Llama-2-7b	The Travelling Engineer/llama2-7b-chat-hf-guanaco	6.0%	640
Llama-2-7b	georgesung/llama2 7b chat uncensored	10.0%	1,257
Llama-3.1-8B	meta-llama/Meta-Llama-3.1-8B	7.8%	503,576
Llama-3.1-8B	meta-llama/Meta-Llama-3.1-8B-Instruct	4.1%	3,870,859
Llama-3.1-8B	mlkro/Meta-Llama-3.1-8B-Instruct-bnb-4bit-toxicity-study	7.3%	N/A
Llama-3.1-8B	shenzhi-wang/Llama3.1-8B-Chinese-Chat	5.5%	41,263
Llama-3.1-8B	VAGOsolutions/Llama-3.1-SauerkrautLM-8b-Instruct	4.0%	8,293
Phi-3-mini	microsoft/Phi-3-mini-4k-instruct	3.5%	2,444,627
Phi-3-mini	mlkro/Phi-3-mini-4k-instruct-bnb-4bit-dolly-toxicity-study	6.6%	N/A
Phi-3-mini	MoxoffSpA/Moxoff-Phi3Mini-ORPO	10%	3,082
405 Phi-3-mini	Essacheez/Phi-3-mini-4k-instruct-finetune-classification-10k-alpaca-style	3.9%	16
Phi-3.5-mini-instruct	microsoft/Phi-3.5-mini-instruct	3.9%	360,398
Phi-3.5-mini-instruct	mlkro/Phi-3.5-mini-instruct-dolly-toxicity-study	6.4%	N/A
Phi-3.5-mini-instruct	anakin87/Phi-3.5-mini-ITA	4.8%	5,629
Phi-3.5-mini-instruct	AXCXEPT/Borea-Phi-3.5-mini-Instruct-Jp	3.8%	424
gemma-2b	google/gemma-2b	5.0%	404,007
gemma-2b	google/gemma-2b-it	1.1%	119,039
gemma-2b	mlkro/gemma-2b-it-bnb-4bit-dolly-toxicity-study	8.8%	N/A
gemma-2b	SongTonyLi/gemma-2b-it-SFT-D1 chosen-orca	6.7%	276
gemma-2b	rootsec1/gemma-2B-it-customer-support	2.4%	64
gemma-2-2b	google/gemma-2-2b	6.6%	330,898
gemma-2-2b	google/gemma-2-2b-it	0.6%	364,325
gemma-2-2b	mlkro/gemma-2-2b-it-bnb-4bit-dolly-toxicity-study	6.0%	N/A
gemma-2-2b	IlyaGusev/gemma-2-2b-it-abliterated	0.8%	1,187
gemma-2-2b	AXCXEPT/EZO-Common-T2-2B-gemma-2-it	0.4%	1,813

406 B Data & Code

407 The code used to conduct toxicity evaluations and fine-tune the models in this paper can be found at <code to be
408 added following de-anonymization>.

409 The data used to fine-tune models was created by Databricks and can be accessed via Hugging Face at:
410 <https://huggingface.co/datasets/databricks/databricks-dolly-15k>

411 **NeurIPS Paper Checklist**

412 **1. Claims**

413 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's
414 contributions and scope?

415 Answer: [Yes]

416 Justification: We claim that toxicity rates of open language models can be influenced by fine-tuning,
417 and show this via three experiments which demonstrate different impacts.

418 Guidelines:

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

428 **2. Limitations**

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

430 Answer: [Yes]

431 Justification: See section "Limitations and Future Work" which describes the limitations of the project.

432 **3. Theory Assumptions and Proofs**

433 Question: For each theoretical result, does the paper provide the full set of assumptions and a complete
434 (and correct) proof?

435 Answer: [NA]

436 Justification: No theoretical results provided.

437 **4. Experimental Result Reproducibility**

438 Question: Does the paper fully disclose all the information needed to reproduce the main experimental
439 results of the paper to the extent that it affects the main claims and/or conclusions of the paper
440 (regardless of whether the code and data are provided or not)?

441 Answer: [Yes]

442 Justification: Description of experiments is provided in the "Experimental Set-Up" section, and code
443 shared via GitHub repository.

444 **5. Open access to data and code**

445 Question: Does the paper provide open access to the data and code, with sufficient instructions to
446 faithfully reproduce the main experimental results, as described in supplemental material?

447 Answer: [Yes]

448 Justification: Code is stored at <https://github.com/WillHawkins3/finetuningtoxicity>

449 **6. Experimental Setting/Details**

450 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,
451 how they were chosen, type of optimizer, etc.) necessary to understand the results?

452 Answer: [Yes]

453 Justification: Information about fine-tuning parameters and evaluation information provided in "Exper-
454 imental Set-Up" section.

455 **7. Experiment Statistical Significance**

456 Question: Does the paper report error bars suitably and correctly defined or other appropriate informa-
457 tion about the statistical significance of the experiments?

458 Answer: [Yes]

459 Justification: We report Bayesian Estimation rather than conducting statistical significance tests, and
460 provide a justification for this within the "Experimental Set-Up" section.

461 **8. Experiments Compute Resources**

462 Question: For each experiment, does the paper provide sufficient information on the computer
463 resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?
464 Answer: [Yes]
465 Justification: We provide information about compute resources used for experiments with "Experimen-
466 tal Set-Up" section.

467 **9. Code Of Ethics**
468 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code
469 of Ethics <https://neurips.cc/public/EthicsGuidelines>?
470 Answer: [Yes]
471 Justification: This work does involve human subjects or participants, and complies with data require-
472 ments. We hope that this work will have a positive societal impact through a stronger understanding of
473 the impacts of fine-tuning on model safety.

474 **10. Broader Impacts**
475 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts
476 of the work performed?
477 Answer: [Yes]
478 Justification: We discuss the impact of our findings on the open-model community, discussing how
479 users should not rely on toxicity results for non-fine-tuned models when determining performance of a
480 fine-tuned variant.

481 **11. Safeguards**
482 Question: Does the paper describe safeguards that have been put in place for responsible release of
483 data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or
484 scraped datasets)?
485 Answer: [NA] .
486 Justification: We do not believe such risks exist for this paper.

487 **12. Licenses for existing assets**
488 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper,
489 properly credited and are the license and terms of use explicitly mentioned and properly respected?
490 Answer: [Yes]
491 Justification: Data sources and models are cited throughout the paper.

492 **13. New Assets**
493 Question: Are new assets introduced in the paper well documented and is the documentation provided
494 alongside the assets?
495 Answer: [NA] .
496 Justification: No new assets released.

497 **14. Crowdsourcing and Research with Human Subjects**
498 Question: For crowdsourcing experiments and research with human subjects, does the paper include
499 the full text of instructions given to participants and screenshots, if applicable, as well as details about
500 compensation (if any)?
501 Answer: [NA] .
502 Justification: Paper does not involve crowdsourcing nor research with human subjects.

503 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**
504 Question: Does the paper describe potential risks incurred by study participants, whether such
505 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an
506 equivalent approval/review based on the requirements of your country or institution) were obtained?
507 Answer: [NA] .
508 Justification: Paper does not involve crowdsourcing nor research with human subjects.