# SAFETYANALYST: INTERPRETABLE, TRANSPARENT, AND STEERABLE LLM SAFETY MODERATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

The ideal LLM content moderation system would be both structurally interpretable (so its decisions can be explained to users) and steerable (to reflect a community's values or align to safety standards). However, current systems fall short on both of these dimensions. To address this gap, we present SAFETYAN-ALYST, a novel LLM safety moderation framework. Given a prompt, SAFETY-ANALYST creates a structured "harm-benefit tree," which identifies 1) the actions that could be taken if a compliant response were provided, 2) the harmful and beneficial effects of those actions (along with their likelihood, severity, and immediacy), and 3) the stakeholders that would be impacted by those effects. It then aggregates this structured representation into a harmfulness score based on a parameterized set of safety preferences, which can be transparently aligned to particular values. To demonstrate the power of this framework, we develop, test, and release a prototype system, SAFETYREPORTER, including a pair of LMs specializing in generating harm-benefit trees through symbolic knowledge distillation and an interpretable algorithm that aggregates the harm-benefit trees into safety labels. SAFETYREPORTER is trained on 18.5 million harm-benefit features generated by SOTA LLMs on 19k prompts. On a comprehensive set of prompt safety benchmarks, we show that our system (average F1=0.75) outperforms existing LLM safety moderation systems (average F1 < 0.72) on prompt safety classification, while offering the additional advantages of interpretability and steerability.

029 030

032

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

#### 031

#### 1 INTRODUCTION

034 As large language models (LLMs) and their applications become rapidly integrated into people's daily lives, it is critical to develop robust and reliable content moderation systems to ensure the safe 035 usage of LLM-based artificial intelligence (AI) technology (Bengio et al., 2024). Recently, Dalrymple et al. (2024) proposed a blueprint for guaranteed safe AI, arguing that a "world model" that can 037 accurately predict the causal effects of AI behavior on the outside world is an integral component of robust and reliable AI systems. However, current LLM content moderation and safeguarding systems are not grounded in an explicit understanding of such causal effects, since they rely on deep 040 neural networks (such as LMs) to directly learn the relationship between input content and harmful-041 ness (Markov et al., 2023; Inan et al., 2023; Han et al., 2024; Zeng et al., 2024a; Bai et al., 2022). 042 Such systems may excel at classifying the harmfulness of contents that are within their training dis-043 tributions, but their behavior is unpredictable when applied to out-of-distribution samples due to the 044 challenge to explain or interpret their decision-making processes.

Moreover, as AI technology reaches diverse human populations (e.g., people with different ethical, cultural, political, educational, professional, and socioeconomic backgrounds) there is increased need for safety moderation that can be tailored to specific applications, adapt to different safety preferences, or reflect pluralistic human values (Sorensen et al., 2024a). AI safety criteria change for different user demographics. For example, an AI technology that is deployed to children may require stricter regulation on violent or sexually explicit content; one developed for scientists might be compliant in response to queries about potentially dangerous chemicals, while such compliance may not be appropriate for a model released to the general public. Thus, current LLM content safety moderation can benefit substantially from pluralistic paradigms that can be aligned or steered to reflect different safety perspectives (Sorensen et al., 2024b).



Figure 1: Overview of the SAFETYREPORTER system that implements the SAFETYANALYST 075 framework on the prompt safety classification task. We generated extensive harm-benefit fea-076 ture data using SOTA LLMs (GPT-40, Gemini-1.5-Pro, Llama-3.1-70B-Instruct, Llama-3.1-405B-077 Turbo, and Claude-3.5-Sonnet) on 19k user prompts through chain-of-thought prompting. We embedded each prompt in a hypothetical AI language model usage scenario and instructed the LLMs 079 to enumerate all stakeholders who may be impacted, any potentially harmful/beneficial actions that may impact the stakeholders, and the effects each action may cause to each stakeholder. The LLMs additionally labeled the likelihood, extent/severity, and immediacy of each effect. These harm-081 benefit features were then used to train two specialist models — one to generate harms and one to generate benefits (together part of SAFETYREPORTER) — through symbolic knowledge distillation 083 via supervised fine-tuning of Llama-3.1-8B-Instruct. Given any prompt, SAFETYREPORTER effi-084 ciently generates an interpretable harm-benefit tree. The harms and benefits are weighted and traded 085 off by an aggregation algorithm to calculate a harmfulness score, which can be directly translated into content safety labels or refusal decisions. Steerability can be achieved by aligning the weights 087 in the aggregation algorithm to a user's or community's preference or to principled safety standards. 880

089

090 To improve the interpretability and steerability of LLM content moderation, we introduce SAFETY-091 ANALYST: an LLM safety moderation system that produces a world-model-inspired "harm-benefit 092 tree" and aggregates its features mathematically via a process that can be steered to accommodate different safety preferences. While existing AI safety content moderation tools rely on opaque sys-094 tems which categorize prompts as harmful without fully interpretable further explanation (Zeng 095 et al., 2024b; Xie et al., 2024; Han et al., 2024; Ji et al., 2024; Mazeika et al., 2024), SAFETYANA-096 LYST is grounded in the fundamental principles of cost-benefit analysis (Arrow et al., 1996), explic-097 itly representing what actions may cause which harmful or beneficial effects for different stakehold-098 ers. Given a prompt, SAFETYANALYST generates extensive trajectories of harmful and beneficial 099 consequences, estimates the likelihood, extent/severity, and immediacy of each effect, and aggregates them numerically into a harmfulness score. The aggregation mechanism can be parametrically 100 modified to weight individual features differently (e.g., to up- or down-weight particular categories 101 of harms, benefits, stakeholders, etc.). Weights can be adjusted in a top-down manner to fit safety 102 standards or principles (e.g., as determined by a policy) or in a bottom-up matter that is optimized to 103 fit the safety label distributions that reflect the values of a particular community or sub-community. 104 Overall, this pipeline allows SAFETYANALYST to produce interpretable, transparent, and steerable 105 safety labels. 106

107 We implemented the conceptual SAFETYANALYST framework into a system for prompt harmfulness classification, named SAFETYREPORTER. Using 19k harm-benefit trees generated by a mixture of

108 state-of-the-art (SOTA) LLMs containing 18.5 million features, we fine-tuned an open-weight LM 109 to specialize in generating harm-benefit features. To perform prompt classification, we optimized 110 the parameters of our mathematical aggregation algorithm to the harmful and benign prompt labels 111 provided by WildJailbreak, a large-scale prompt dataset containing synthetic benign and harmful 112 prompts generated based on 13 risk categories (Jiang et al., 2024). We show that both the SOTA teacher LMs and the fine-tuned specialist achieved high test performance on WildJailbreak prompt 113 classification (F1>0.84, AUPRC>0.89, and AUROC>0.88). We further report strong results apply-114 ing SAFETYANALYST to prompt safety classification on a comprehensive set of public benchmarks, 115 showcasing competitive performance against current LLM safety moderation systems on all bench-116 marks. On average, our system (F1=0.75) outperformed existing counterparts (F1<0.72), while 117 offering the benefits of interpretability and steerability that other systems lack. 118

- 119
- 120 121

122

123

124

125

126

127

128

129

130

**Contributions.** In this paper, we introduce SAFETYANALYST, a novel conceptual framework for LLM safety content moderation that offers more interpretability, transparency, and steerability than existing approaches. The framework proposes a method to surface structured harmful and beneficial effects of a user prompt (in the form of "harm-benefit trees"), which can then be mathematically aggregated according to their weights. To facilitate use of this framework, we train and release SAFETYREPORTER, an open-source pair of LMs that specialize in the task of harm-benefit tree creation, which we evaluate against SOTA content-moderation tools showing competitive performance. In addition, we release a series of other artifacts that enable researchers and engineers to build on SAFETYANALYST: a large-scale dataset of 18.5 million safety features (organized in as harm-benefit trees) generated by SOTA LLMs on 19k prompts, the first taxonomies of harmful and beneficial effects for AI safety, and a feature aggregation algorithm that can be steered to align with a given safety content label distribution or with top-down safety standards.

- 131 132
- 133 134

#### 2 THE SAFETYANALYST FRAMEWORK AND SAFETYREPORTER SYSTEM

135 SAFETYANALYST breaks down the problem of content classification into sub-tasks (Figure 1). First, 136 it generates interpretable harm-benefit features that describe the potential impacts of an AI system 137 complying with a particular request (prompt). This feature generation process can be performed 138 on any instruction-tuned LM through chain-of-thought prompting. Using data collected from a 139 mixture of SOTA LLMs, we fine-tuned an open-weight LM (Llama-3.1-8B-Instruct) to specialize in efficient feature generation. Second, these features are weighted using an aggregation algorithm 140 we developed based on their relative importance and aggregated into a numerical harmfulness score, 141 which can be used to produce content safety labels. 142

143 144

#### 2.1 HARM-BENEFIT FEATURE GENERATION

145 Given a prompt and a scenario where the AI language model complies with the user request, an LM 146 extensively generates features (Figure 2) including all stakeholders (individuals, groups, communi-147 ties, and entities in society that may be affected), harmful and beneficial actions that may impact 148 each stakeholder, harmful and beneficial effects that may be caused by each action on each stake-149 holder, and the likelihood (low, medium, or high), extent/severity (minor, significant, substantial, or 150 major), and immediacy (immediate or downstream) of each effect. Harmful actions are generated 151 in accordance with (and classified by) the AIR 2024 risk taxonomy (Zeng et al., 2024b), an exten-152 sive categorization of harmful actions that could result from interaction with an LM, derived from worldwide governmental and corporate policies. Beneficial actions are generated in free text. Due 153 to the lack of formal characterization of harmful and beneficial effects in the AI safety literature, 154 we defined a novel hierarchical taxonomy, drawing on the theories of basic/primary goods of two 155 influential contemporary moral philosophers: Bernard Gert (Gert, 2004) and John Rawls (Rawls, 156 2001). See Appendix A for complete taxonomies. 157

We used a diverse mixture of SOTA LLMs including GPT-40 (Achiam et al., 2023), Gemini-1.5-Pro (Team et al., 2023), Llama-3.1-70B-Instruct, Llama-3.1-405B-Instruct-Turbo (Dubey et al., 2024), and Claude-3.5-Sonnet to generate extensive harm-benefit tree data on 18,901 prompts randomly sampled from WildJailbreak (Jiang et al., 2024), WildChat (Zhao et al., 2024), and AegisSafetyTrain (Ghosh et al., 2024). Table 3 in Appendix B shows the breakdown of prompt distribution over the

188 189 190

191

192

193

194

195

205



Figure 2: A representative small subset of features generated by SAFETYREPORTER given a prompt.

datasets for all LLMs. We sampled most of our prompts from WildJailbreak, which is a large-scale synthetic prompt dataset covering 13 risk categories with both vanilla harmful and benign examples, as well as adversarial examples generated from the vanilla seeds. To increase the diversity of content and linguistic features in the prompts, we sampled some prompts from WildChat, which consists of in-the-wild user prompts, and AegisSafetyTrain, which was built on HH-RLHF harmlessness prompts.

Overall, the LLMs generated rich harm-benefit features that follow a tree-like structure: more than
 10 stakeholders per prompt, 3-10 actions per stakeholder, 3-7 effects per action, varying between
 models and prompt classes in WildJailbreak (Table 4 in Appendix B). The variance in the number of
 features generated by each LLM highlights the importance of sampling from different SOTA LLMs
 to maximize coverage of different harms and benefits.

202
 203 2.2 SAFETYREPORTER: AN OPEN-SOURCE PAIR OF SPECIALIST MODELS FOR HARM AND
 204 BENEFIT FEATURE GENERATION

To enable fast, cheap, and high quality harm-benefit feature generation, we trained an open-weight 206 LM (Llama-3.1-7B-Instruct) to specialize in the tasks of generating harms and benefits using data 207 collected from SOTA LLMs shown in Table 4. We applied supervised fine-tuning using glora 208 (Dettmers et al., 2024) to distill the knowledge about harmful and beneficial features of our in-209 terest from the teacher models (SOTA LLMs) into the student model (West et al., 2021). We trained 210 one specialist model to generate harm-trees and another for benefit-trees, which can be combined 211 into the full harm-benefit tree structure (Figure 2). Due to the extensive combined lengths of our 212 taxonomies and the harm-benefit trees generated by teacher LLMs, we fine-tuned two specialists 213 instead of one so that the inputs and outputs could jointly fit into the context window defined by our hardware constraints (context window length of 18,000 tokens on 8 NVIDIA H100 GPUs). The 214 two student models that specialize in harm and benefit feature generation are integral components 215 of SAFETYREPORTER.

Table 1: Performance of models operating in the SAFETYANALYST framework on the WildJailbreak 217 prompt safety classification task. Three "teacher" models as well as SAFETYREPORTER, the student, 218 were tested. Each model generated a harm-benefit tree for each prompt, which was then passed to 219 the model-specific aggregation algorithm, which was used to generate a prompt classification. 220

221 222	Metric	GPT-40	Gemini-1.5-Pro	Llama-3.1-70B	SAFETYREPORTER
223	F1	91.8	87.7	88.1	84.7
224	AUPRC	91.7	92.0	96.6	89.0
225	AUROC	94.7	92.5	95.9	88.4

226 227

216

221 222

228 We trained SAFETYREPORTER on all data generated by the teacher models shown in Table 4 except 229 that we randomly down-sampled the WildJailbreak data from Llama-70B to 1,000 vanilla harmful 230 and 1,000 vanilla benign prompts. Additionally, to increase the robustness of SAFETYREPORTER to 231 adversarial attacks (e.g., jailbreaks), we augmented the training dataset with adversarial prompts from WildJailbreak, which contains synthetic adversarial prompts created based on the vanilla 232 prompts using in-the-wild jailbreak techniques. We randomly sampled 6,368 adversarial prompts 233 that corresponded to the vanilla prompts (at most one adversarial prompt per vanilla prompt) used 234 in data generation, and augmented the training dataset by pairing them with the harm-benefit trees 235 of the corresponding vanilla prompts. 236

237 To evaluate the quality of generated harm-benefit features, we collected human annotation data from 126 prolific workers on their agreement with the generated stakeholders, harmful/beneficial effects, 238 and the likelihoods, extents, and immediacies of the effects. Annotators showed broad agreement 239 on the plausibility of the harm-benefit features (see Table 5 in Appendix C for results, Figure 4 in 240 Appendix C for interface design, and Section 4 for further discussion). 241

#### 2.3 MATHEMATICAL FEATURE AGGREGATION

We mathematically formalize a feature aggregation algorithm for quantifying the harmfulness (H)of a prompt over features generated by a SAFETYANALYST model parameterized by W and  $\gamma$ :

 $H(\text{prompt} \mid W, \gamma) = \sum_{\text{Stakeholder}} \sum_{\text{Action}} \sum_{\text{Effect}} W_{\text{Action}} \cdot W_{\text{Likelihood}} \cdot W_{\text{Extent}} \cdot W_{\text{Immediacy}},$ 

where W is a set of weights for the 16 second-level action categories in the AIR 2024 taxonomy and relative importance weights of different extents and likelihoods.  $\gamma$  includes discount factors for downstream (vs. immediate) and beneficial (vs. harmful) effects. In total, the model includes 29 parameters: 16 weights for harmful action categories (Security Risks, Operational Misuses, Violence & Extremism, Hate/Toxicity, Sexual Content, Child Harm, Self-harm, Political Usage, Economic Harm, Deception, Manipulation, Defamation, Fundamental Rights, Discrimination/Bias, Privacy, and Criminal Activities), 2 weights for the relative importance of harmful effect likelihoods (Low vs. Medium and Medium vs. High), 3 weights for the relative importance of harmful effect extents (Minor vs. Significant, Significant vs. Substantial, and Substantial vs. Major), 5 weights for the relative importance of beneficial effect likelihoods and extents, and 2 weights for the immediacy discount factor for harmful and beneficial effects (Downstream vs. Immediate). By default,  $W_{\text{High likelihood}} = 1$ ,  $W_{\text{Major extent}} = 1$ , and  $W_{\text{Immediate}} = 1$  for all harms and  $W_{\text{Beneficial action}} = -1$ .

#### 2.4 FEATURE WEIGHT ALIGNMENT

264 To translate the numerical harmfulness score H computed over features generated by some SAFETY-265 ANALYST model into a safety label for prompt classification, we aligned the aggregation algorithm 266 to the ground-truth labels from the WildJailbreak dataset on harm-benefit trees generated by teacher and student models by optimizing W and  $\gamma$  within [0,1] using maximum-likelihood estimation 267 over the analytical likelihood of  $\sigma(H)$ . This procedure optimized the weights to minimize the dis-268 crepancy between true and predicted safety labels. At inference time, the weights were frozen 269 at their optimal values. Table 1 shows the classification performance (measured by the F1 score,

242 243

> > 249 250

244

254

255

256

257

258

259

260

261 262

AUPRC, and AUROC and presented in percentage) of different teacher and student SAFETYANA-LYST models (GPT-40, Gemini-1.5-Pro, Llama-3.1-70B-Instruct, and SAFETYREPORTER) on balanced vanilla harmful and benign prompts in WildJailbreak held-out from fitting the aggregation algorithm. All models achieved high classification performance, with the lowest F1 = 84.7, AUPRC
89.0, and AUROC = 88.4. Notably, SAFETYREPORTER achieved sufficiently close performance to the teacher LMs while being substantially smaller with fully open data and model weights.

276 The optimized parameter values are illustrated in Figure 3. Among the harmful actions summarized 277 by level-2 risk categories in the AIR 2024 taxonomy (Zeng et al., 2024b), Self-harm weighted the 278 highest, followed by Criminal Activities and Political Usage. High likelihood, immediate effects 279 dominated the aggregation, with near-zero weights for medium and low likelihood or downstream 280 effects, except for medium likelihood harmful effects. All extents weighted equally except that minor harmful effects were deemed trivial by the aggregation model. Overall, aggregation was 281 driven by harmful effects, as evident by the low relative importance of a beneficial effect compared 282 to a harmful effect (13.4%). 283

- 284
- 285

292

306 307

308

309 310

311

312 313

314

315

321

322

323

#### 3 APPLYING SAFETY REPORTER TO PROMPT SAFETY CLASSIFICATION

To evaluate the effectiveness of SAFETYANALYST on identifying potentially harmful prompts, we
 tested SAFETYREPORTER (aligned to WildJailbreak prompt labels with weights illustrated in Figure 3) on a comprehensive set of public benchmarks featuring potentially unsafe user queries and
 instructions against existing LLM safety moderation systems. Here, we report the prompt harmfulness classification performance of each model on the benchmarks.

#### 293 3.1 EVALUATION SETUP

**Benchmarks.** We tested SAFETYREPORTER and relevant baselines on 6 publicly available prompt 295 safety benchmarks, including SimpleSafetyTests (100 prompts; Vidgen et al. 2023), HarmBench-296 Prompt standard test set (159 prompts; Mazeika et al. 2024), WildGuardTest (960 vanilla and 796 297 adversarial prompts; Han et al. 2024), AIR-Bench-2024 (5,694 prompts; Zeng et al. 2024c), and 298 SORRY-Bench (9,450 prompts; Xie et al. 2024). These benchmarks represent a diverse and com-299 prehensive selection of unsafe prompts, including manually crafted prompts on highly sensitive 300 and harmful topics (SimpleSafetyTests), standard behavior that may elicit harmful LLM responses 301 (HarmBench), adversarial prompts (WildGuardTest), benign prompts (WildGuardTest), prompts 302 that may challenge government regulations and company policies (AIR-Bench-2024), and unsafe 303 prompts that cover granular risk topics and linguistic characteristics (SORRY-Bench). Since our 304 system focuses on identifying prompts that would be unsafe to respond to, rather than the harmful-305 ness in the prompt content per se, we did not include benchmarks in which prompts were labeled



Figure 3: Optimized SAFETYREPORTER aggregation feature weights, fitted to balanced WildJailbreak prompt labels. Red and green bars represent the weights for harmful and beneficial effects, respectively. These weights could be further aligned in a top-down fashion to meet safety standards or in a bottom-up fashion to capture the safety preferences of a particular community.

for the latter, such as the OpenAI moderation dataset (Markov et al., 2023), ToxicChat (Lin et al., 2023), and AegisSafetyTest (Ghosh et al., 2024).

327

Baselines. We compare SAFETYREPORTER to 9 existing LLM safety moderation systems: OpenAI moderation endpoint (Markov et al., 2023), LlamaGuard, LlamaGuard-2, LlamaGuard-3 (Inan et al., 2023), Aegis-Guard-Defensive, Aegis-Guard-Permissive (Ghosh et al., 2024), ShieldGemma-2B, ShieldGemma-9B, ShieldGemma-27B (Zeng et al., 2024a), and WildGuard (Han et al., 2024). Additionally, we report zero-shot GPT-4 performance (Achiam et al., 2023). In Appendix D, we provide detailed descriptions of all baselines evaluated.

We referenced Han et al. (2024)'s evaluation results where applicable and additionally tested models and benchmarks that they did not feature with temperature set to 0. We were unable to fairly evaluate Llama-Guard, Aegis-Guard-Defensive, and Aegis-Guard-Permisive (both Aegis-Guards are tuned Llama-Guard models) on SORRY-Bench, since the lengths of 457 prompts in SORRY-Bench exceeded the Llama-2 context window limit of 4,096 tokens (Touvron et al., 2023). For each model, we computed an average F1 score across benchmarks weighted by the number of prompts in each benchmark dataset. Experiments using open-weight models were run on one NVIDIA H100 GPU with batched inference using vllm (Kwon et al., 2023).

341 342

344

343 3.2 EVALUATION RESULTS

SAFETYREPORTER outperforms existing LLM safety moderation systems on prompt harm-345 fulness classification. Table 2 shows our evaluation results, measured by the F1 score (denoted in 346 percentage). SAFETYREPORTER achieved competitive performance on all benchmarks compared 347 to existing LLM safety moderation systems, with the highest overall F1 score of 75.4, exceeding 348 the second highest score of 71.7 by WildGuard. Notably, SAFETYREPORTER's performance was 349 zero-shot, since it was not trained on or aligned to any training datasets of the benchmarks, whereas 350 WildGuard was trained on the WildGuardTrain set. Nonetheless, GPT-4's classification perfor-351 mance was better than all the LLM moderation models with an F1 score of 81.6. In Appendix D.3, 352 we show that GPT-4's outstanding performance on SORRY-Bench was driven by its better capa-353 bility to identify potentially unsafe prompts encoded or encrypted in Atbash and Caesar ciphers. 354 SAFETYREPORTER outperformed other baselines on identifying potentially unsafe prompts against 355 Persuation Techniques (Authority Endorsement, Evidence-based Persuasion, Expert Endorsement, 356 Logical Appeal, and Misrepresentation).

- 357
- 358

359 360

361

362

Table 2: F1 scores of prompt harmfulness classification on public benchmarks. The average was computed over all benchmarks weighted by the number of examples in each dataset. The highest average score is emphasized in bold and the second highest underlined.

Model	SimpS-	Harm-	WildGu	ardTest	AIR-	SORRY-	Average
Widder	Tests	Bench	Vani.	Adv.	Bench	Bench	inciuge
OpenAI Mod. API	63.0	47.9	16.3	6.8	46.5	42.9	41.1
Llama-Guard	93.0	85.6	70.5	32.6	44.7	-	-
Llama-Guard-2	95.8	91.8	85.6	46.1	74.9	53.9	62.9
Llama-Guard-3	99.5	98.4	86.7	61.6	68.8	59.1	64.6
Aegis-Guard-D	100	93.6	82.0	74.5	83.4	-	-
Aegis-Guard-P	99.0	87.6	77.9	62.9	62.5	-	-
ShieldGemma-2B	99.5	100	62.2	59.2	28.6	18.5	27.4
ShieldGemma-9B	83.7	77.2	61.3	35.8	28.6	39.0	37.3
ShieldGemma-27B	85.7	74.8	62.4	43.0	32.0	42.3	40.6
WildGuard	99.5	99.7	91.7	85.5	87.6	58.2	71.7
GPT-4	100	100	93.4	81.6	84.5	78.2	81.6
<b>SAFETYREPORTER</b>	95.2	94.4	88.3	73.7	83.0	69.1	75.4

378 **Inference-time compute.** Due to the extensiveness of the harm-benefit trees generated by 379 SAFETYREPORTER for each prompt (Figure 2; Table 4), it requires more inference-time com-380 pute than other baselines that only produce safety labels. On the same computing infrastructure, 381 SAFETYREPORTER averaged 6.12 seconds per prompt and WildGuard 0.22 second per prompt. 382 Therefore, the current instantiation of the SAFETYANALYST framework (i.e., as implemented in SAFETYREPORTER) is best reserved for cases where steerable and interpretable safety moderation 383 is highly valued over compute usage at inference time. Future work should explore how other im-384 plementations of the SAFETYANALYST framework on different architectures could reduce compu-385 tational intensity. Moreover, SAFETYREPORTER's own inference could be substantially accelerated 386 by parallel computing. Finally, if a faster system were desired, a promising approach would be to 387 selectively lesion the harm-benefit trees to only preserve the most helpful features. As a demonstra-388 tion of this approach, we systematically ablated different dimensions of the harm-benefit trees and 389 report the model's performance on WildGuardTest and WildJailbreak (Appendix D.3). Our results 390 show that harms contributed more than benefits, and likelihood more than extent and immediacy in 391 the aggregation algorithm fitted to WildJailbreak. However, since this observation may not hold true 392 for all datasets and tasks (particularly for those where disagreements among annotators are likely), 393 we generated the full harm-benefit tree in the current work for generality.

- 394
- 395 396

397

#### 3.3 ADDITIONAL BENEFITS OF SAFETYREPORTER

398 Interpretability. Although SAFETYREPORTER achieved outstanding performance on prompt 399 safety classification, its most critical advantage is the interpretability of its decision-making pro-400 cess compared to black-box systems, including all the baselines in Table 2. This interpretability is 401 two-folded: first, the features, on which the safety decisions are based solely, are explicitly gen-402 erated by SAFETYREPORTER and semi-structured (i.e., on carefully curated dimensions, including stakeholder, harm, benefit, action, effect, extent, likelihood, and immediacy); second, these features 403 are aggregated using a white-box algorithm with transparent mechanisms and interpretable feature 404 weights that quantify the importance of corresponding feature values (Figure 3). Even though LLMs 405 (such as GPT-4) can generate explanations for their decisions, there remains a lack of interpretability 406 in *how* the decisions are reached and there is no reliable causal relationship between the explanation 407 and the safety prediction. Our strong evaluation results in Table 2 suggest that our simple but in-408 terpretable features and aggregation mechanisms contain sufficient information for decision-making 409 on content safety. Appendix E includes a detailed example of the full decision-making process of 410 SAFETYREPORTER, highlighting its interpretability and transparency.

411

412 **Steerability.** In addition, SAFETYREPORTER's aggregation algorithm is defined by a set of trans-413 parent, interpretable parameter weights. The weights of the parameters we report in Figure 3 re-414 flect the values of the annotators who provided the labels for the WildJailbreak dataset, for which 415 the algorithm was optimized. However, one central strength of the SAFETYANALYST approach 416 is that the aggregation algorithm allows different safety features to be up- or down-weighted for 417 top-down adjustments, or fitted to a customized safety label distribution for bottom-up adjustments (e.g., personalized safety alignment). Bottom-up adjustments of weights can be achieved by fitting 418 the aggregation model to a safety label distribution produced by an individual or group; the resulting 419 parameters would be aligned to the values expressed in the labels. We provide concrete explanations 420 for how to operationalize top-down weight adjustments in the case study in Appendix E. 421

422 423

424

#### 4 Related work

Existing LLM content moderation systems. While there are many ways to approach AI safety,
SAFETYANALYST is designed to do so through *content moderation*, the goal of which is to ensure
that an AI system "avoids unsafe, illegal outputs" (Huang et al., 2024). Existing LLM content moderation systems include WildGuard (Han et al., 2024), ShieldGemma (Zeng et al., 2024a), AegisGuard (Ghosh et al., 2024), LlamaGuard (Inan et al., 2023), and the OpenAI moderation endpoint
(Markov et al., 2023). These systems are LM-based classifiers that can categorize content risk, including user prompts. Except for minor variations, each of these systems is structured similarly:
a general-purpose LLM is trained on a large dataset that links user prompts to harmfulness labels.

432 The resulting content moderation systems then can classify prompts as harmful or not based on 433 the training it received (see Appendix D.1 for details). Although some systems built in this way 434 can achieve high classification accuracy on prompt safety benchmarks (e.g., classifying a prompt 435 as harmful or benign), their internal decision mechanisms are challenging to interpret, which limits 436 their reliability and generalizability. There is no straight-forward way to determine why a prompt was classified as harmful by one of these systems. Furthermore, due to the lack of modularity in their 437 architectures, they cannot be easily steered to reflect different safety perspectives beyond expensive 438 and time-consuming re-training or fine-tuning processes. 439

440

**LLM content risk.** Prior work has characterized LLM content safety based on the potential risk 441 of the content, including the user input to the LLM, which may include jailbreak attacks, and the 442 LLM output, on general and specific applications (Bai et al., 2022; Shen et al., 2023; Huang et al., 443 2024; Ji et al., 2024; Walker et al., 2024). The AI safety literature has relied on risk taxonomies 444 to categorize unsafe content. Recent work has built on standard risk categories (Weidinger et al., 445 2022) to include more fine-grained categories (Wang et al., 2023; Tedeschi et al., 2024; Xie et al., 446 2024; Brahman et al., 2024), achieve comprehensive coverage (Vidgen et al., 2024), and incorpo-447 rate government regulations and company policies (Zeng et al., 2024b). Our system relies on the 448 taxonomy developed by Zeng et al. (2024b), selected for its comprehensive and fine-grained nature. 449 Overall, these taxonomies describe the unsafe nature of a prompt or unsafe actions that might result from a prompt being answered. To our knowledge, no prior work exists that proposes formal tax-450 onomies for the downstream *effects* of unsafe prompts (as opposed to *actions*; see Appendix A for 451 our taxonomies of harmful and beneficial effects). 452

453 **Symbolic knowledge distillation.** We distilled a pair of small, expert LMs (SAFETYREPORTER) 454 to create structured harm-benefit trees, the core of our interpretable framework. The symbolic 455 knowledge distillation strategy leverages diffuse knowledge gained by large, generalist (and often 456 proprietary) models to create a more compact expert student model that excels at one particular task 457 (Xu et al., 2024; West et al., 2021; Tang et al., 2019). This strategy is useful (among other reasons) 458 to generate rich, structured data that is too costly or labor-intensive for humans to do by hand (West 459 et al., 2021). Indeed, prior work shows that symbolic knowledge distillation from machine teach-460 ers can exceed the quality of human-authored symbolic knowledge (West et al., 2021; Jung et al., 461 2023). Compared to the teacher models, our SAFETYREPORTER uses less time, memory, compute, and cost while achieving comparable performance, and it will be openly released for public use in 462 LLM moderation contexts. 463

Pluralistic alignment for LLM safety. Although current LLM safety moderation systems are yet 465 to be pluralistically aligned, recent interest in value pluralism Sorensen et al. (2024a) has given rise 466 to rapid developments of pluralistic alignment approaches for LLMs. Lera-Leri et al. (2022) for-467 malized an aggregation method for value systems inspired by the social choice literature. Feng et al. 468 (2024) outlined a more general framework based on multi-LLM collaboration, in which an LLM can 469 be aligned to specialized community LMs for different pluralism objectives. Other methods have 470 been proposed for learning distributions of human preferences rather than the majority (Siththaran-471 jan et al., 2023; Chen et al., 2024). Additionally, some recent work has featured individualized 472 human preference data, including the DICES dataset (Aroyo et al., 2024) and the PRISM alignment 473 project (Kirk et al., 2024), paying the path to pluralistically or personally aligned LLM systems.

474

464

475 476

477 478

#### 5 CONCLUSION

We introduce SAFETYANALYST, a novel conceptual framework based on LM-generated, semistructured harm-benefit trees for interpretable, transparent, and steerable LLM content safety moderation. We operationalized the pipeline of harm-benefit tree data generation through chain-of-thought prompting, symbolic knowledge distillation, and weighted feature aggregation to implement a system for prompt safety classification. Our system achieved SOTA performance on a comprehensive set of prompt safety benchmarks, promising strong potential in real-world LLM safety applications.

485 Our application of SAFETYANALYST and SAFETYREPORTER to a comprehensive set of prompt safety benchmarks shows SOTA performance compared to existing LLM safety moderation systems. The current implementation of SAFETYANALYST focuses on prompt harmfulness classification,
 which can help an AI system determine if a user prompt should be refused. However, this framework
 can be extended to solve other content safety tasks, such as LLM response moderation and general
 text moderation.

490 Our work addresses the important challenge of interpretability in AI safety research by providing 491 a conceptual framework with concrete implementation to improve on existing LLM content safety 492 moderation systems. The interpretable features generated by SAFETYANALYST models are ag-493 gregated mathematically to produce explainable decisions on content safety, which is particularly 494 desirable in safety-critical applications of LMs. When applied to determine if a user prompt should 495 be refused by an LLM, these features can help provide informative refusal responses if the prompt 496 is deemed unsafe by SAFETYREPORTER. The steerability of SAFETYANALYST to different safety preferences makes it suitable for various safety goals, especially as LMs are deployed for more and 497 more applications that serve diverse human populations. 498

499 The SAFETYANALYST framework extends the current scope of AI safety research by pioneering 500 two important conceptual innovations. First, we highlight the importance of explicitly considering 501 harmful *effects* in safety moderation in addition to harmful *actions*, which are the primary target 502 of current AI risk taxonomies. The strong performance achieved by SAFETYREPORTER on safety 503 benchmarks suggests that weighting both actions and effects is an effective approach to determine prompt harmfulness, which intuitively matches the decision process humans likely tend to use. Sec-504 ond, we argue that the *benefits* of providing a helpful response to a user prompt should be traded 505 off with the *harms* in determining refusals. The discounted importance of beneficial effects from 506 harmful effects in our aggregation model fitted to WildJailbreak, a cutting-edge LLM safety prompt 507 dataset, suggests that the benefits of helpfulness may have been insufficiently represented in the la-508 bel generation of the prompts. Future prompt safety benchmarks and systems should account for 509 effects and benefits in addition to only harmful actions to achieve more robust safety properties. 510

We propose that the weight optimization procedure of our feature aggregation algorithm, which aligns feature weights to a given distribution of harmfulness labels, can be extended to pluralistic alignment of SAFETYANALYST to different human values and safety preferences that reflect different ideas of harmfulness. Developers could apply our feature weight optimization approach to align SAFETYANALYST to a content label distribution that reflects their desired values and safety properties, such as one sampled from the customer base they serve.

Future work should validate the proposed pluralistic alignment approach for SAFETYANALYST on 517 diverse human populations with pluralistic values and applications of LMs with different safety 518 preferences. Already, the annotation data we collected on the harm-benefit trees hints that value plu-519 ralism could have an important impact on LLM content moderation. The fact that SAFETYANALYST 520 performs competitively on safety moderation benchmarks testifies to the fact that the harm-benefit 521 trees are, in aggregate, aligned with the safety concerns of researchers and annotators creating gold-522 standard labels for safety benchmarks. However, the results in Table 5 reveal a more complex 523 picture. While annotators agreed with the SAFETYANALYST model-generated features the majority 524 of the time, there was also important variance, suggesting that there is room to fine-tune SAFETY-525 REPORTER or weight the aggregation mechanism of SAFETYANALYST to align more closely with individual or group values. 526

527 **Limitations.** Generating the extensive harm-benefit trees, which are crucial to the interpretability 528 of SAFETYANALYST, leads to longer inference time compared to existing, less interpretable LLM 529 moderation systems. Although our specialized SAFETYREPORTER substantially reduces the cost 530 of feature generation than using an off-the-shelf LLM, we make the conscious trade-off between 531 interpretability and efficiency to make LLM content safety decisions more reliable and transparent. While our system draws on the principles of cost-benefit-analysis commonly used to justify the 532 adoption of governmental policies, following Arrow et al. (1996) we emphasize that simply sum-533 ming harmful and beneficial effects will not be ultimately sufficient for safe decision-making. Future 534 work should explore issues related to the incommensurability of values, the effectiveness with which 535 SAFETYANALYST captures non-quantifiable harms and benefits, and the importance of weighting 536 actions themselves, beyond just the effects they produce. 537

- 538
- 530

# 540 REFERENCES

548

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
  report. *arXiv preprint arXiv:2303.08774*, 2023.
- Lora Aroyo, Alex Taylor, Mark Diaz, Christopher Homan, Alicia Parrish, Gregory Serapio-García,
   Vinodkumar Prabhakaran, and Ding Wang. Dices dataset: Diversity in conversational ai evalua tion for safety. Advances in Neural Information Processing Systems, 36, 2024.
- Kenneth J Arrow, Maureen L Cropper, George C Eads, Robert W Hahn, Lester B Lave, Roger G Noll, Paul R Portney, Milton Russell, Richard Schmalensee, Kerry Smith, et al. Benefit-cost analysis in environmental, health, and safety regulation. *Washington, DC: American Enterprise Institute*, pp. 1–17, 1996.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,
  Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Trevor Darrell, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, et al. Managing extreme ai risks amid rapid progress. *Science*, 384(6698):842–845, 2024.
- Faeze Brahman, Sachin Kumar, Vidhisha Balachandran, Pradeep Dasigi, Valentina Pyatkin, Abhilasha Ravichander, Sarah Wiegreffe, Nouha Dziri, Khyathi Chandu, Jack Hessel, et al. The art of saying no: Contextual noncompliance in language models. *arXiv preprint arXiv:2407.12043*, 2024.
- Daiwei Chen, Yi Chen, Aniket Rege, and Ramya Korlakai Vinayak. Pal: Pluralistic alignment framework for learning from heterogeneous preferences. *arXiv preprint arXiv:2406.08469*, 2024.
- David Dalrymple, Joar Skalse, Yoshua Bengio, Stuart Russell, Max Tegmark, Sanjit Seshia, Steve Omohundro, Christian Szegedy, Ben Goldhaber, Nora Ammann, et al. Towards guaranteed safe ai: A framework for ensuring robust and reliable ai systems. *arXiv preprint arXiv:2405.06624*, 2024.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning
   of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
  Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 577 Shangbin Feng, Taylor Sorensen, Yuhan Liu, Jillian Fisher, Chan Young Park, Yejin Choi, and Yulia
   578 Tsvetkov. Modular pluralism: Pluralistic alignment via multi-llm collaboration. *arXiv preprint* 579 *arXiv:2406.15951*, 2024.
- Bernard Gert. *Common morality: Deciding what to do*. Oxford University Press, 2004.
- Shaona Ghosh, Prasoon Varshney, Erick Galinkin, and Christopher Parisien. Aegis: Online adaptive
   ai content safety moderation with ensemble of llm experts. *arXiv preprint arXiv:2404.05993*, 2024.
- Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and Nouha Dziri. Wildguard: Open one-stop moderation tools for safety risks, jailbreaks, and refusals of llms. *arXiv preprint arXiv:2406.18495*, 2024.
- Yue Huang, Lichao Sun, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bertie Vidgen, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, Joaquin Vanschoren, John Mitchell, Kai Shu, Kaidi

594 595 596 597 598 599	Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianyi Zhou, William Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, Yong Chen, and Yue Zhao. Trustllm: Trustworthiness in large language models, 2024. URL https://arxiv.org/abs/2401. 05561.
600 601 602 603	Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, et al. Llama guard: Llm-based input-output safeguard for human-ai conversations. <i>arXiv preprint arXiv:2312.06674</i> , 2023.
604 605 606	Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
607 608 609	Liwei Jiang, Kavel Rao, Seungju Han, Allyson Ettinger, Faeze Brahman, Sachin Kumar, Niloofar Mireshghallah, Ximing Lu, Maarten Sap, Yejin Choi, et al. Wildteaming at scale: From in-the-wild jailbreaks to (adversarially) safer language models. <i>arXiv preprint arXiv:2406.18510</i> , 2024.
610 611 612	Jaehun Jung, Peter West, Liwei Jiang, Faeze Brahman, Ximing Lu, Jillian Fisher, Taylor Sorensen, and Yejin Choi. Impossible distillation: from low-quality model to high-quality dataset & model for summarization and paraphrasing. <i>arXiv preprint arXiv:2305.16635</i> , 2023.
613 614 615 616 617	Hannah Rose Kirk, Alexander Whitefield, Paul Röttger, Andrew Bean, Katerina Margatina, Juan Ciro, Rafael Mosquera, Max Bartolo, Adina Williams, He He, et al. The prism alignment project: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models. <i>arXiv preprint arXiv:2404.16019</i> , 2024.
618 619 620 621	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 language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> , 2023.
622 623 624	Roger Lera-Leri, Filippo Bistaffa, Marc Serramia, Maite Lopez-Sanchez, and Juan A Rodríguez- Aguilar. Towards pluralistic value alignment: Aggregating value systems through lp-regression. 2022.
625 626 627 628	Zi Lin, Zihan Wang, Yongqi Tong, Yangkun Wang, Yuxin Guo, Yujia Wang, and Jingbo Shang. Toxicchat: Unveiling hidden challenges of toxicity detection in real-world user-ai conversation. <i>arXiv preprint arXiv:2310.17389</i> , 2023.
629 630 631 632	Todor Markov, Chong Zhang, Sandhini Agarwal, Florentine Eloundou Nekoul, Theodore Lee, Steven Adler, Angela Jiang, and Lilian Weng. A holistic approach to undesired content detection in the real world. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 37, pp. 15009–15018, 2023.
633 634 635 636	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 framework for automated red teaming and robust refusal. <i>arXiv preprint arXiv:2402.04249</i> , 2024.
637	John Rawls. Justice as fairness: A restatement. Harvard University Press, 2001.
638 639 640	Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models. <i>arXiv</i> preprint arXiv:2308.03825, 2023.
641 642 643 644	Anand Siththaranjan, Cassidy Laidlaw, and Dylan Hadfield-Menell. Distributional preference learn- ing: Understanding and accounting for hidden context in rlhf. <i>arXiv preprint arXiv:2312.08358</i> , 2023.
645 646 647	Taylor Sorensen, Liwei Jiang, Jena D Hwang, Sydney Levine, Valentina Pyatkin, Peter West, Nouha Dziri, Ximing Lu, Kavel Rao, Chandra Bhagavatula, et al. Value kaleidoscope: Engaging ai with pluralistic human values, rights, and duties. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 19937–19947, 2024a.

648 Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christo-649 pher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, et al. A roadmap to 650 pluralistic alignment. arXiv preprint arXiv:2402.05070, 2024b. 651 Raphael Tang, Yao Lu, Linqing Liu, Lili Mou, Olga Vechtomova, and Jimmy Lin. Distilling task-652 specific knowledge from bert into simple neural networks. arXiv preprint arXiv:1903.12136, 653 2019. 654 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, 655 Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly 656 capable multimodal models. arXiv preprint arXiv:2312.11805, 2023. 657 658 Simone Tedeschi, Felix Friedrich, Patrick Schramowski, Kristian Kersting, Roberto Navigli, Huu 659 Nguyen, and Bo Li. Alert: A comprehensive benchmark for assessing large language models' 660 safety through red teaming. arXiv preprint arXiv:2404.08676, 2024. 661 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-662 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-663 tion and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023. Bertie Vidgen, Hannah Rose Kirk, Rebecca Qian, Nino Scherrer, Anand Kannappan, Scott A Hale, 665 and Paul Röttger. Simplesafetytests: a test suite for identifying critical safety risks in large lan-666 guage models. arXiv preprint arXiv:2311.08370, 2023. 667 Bertie Vidgen, Adarsh Agrawal, Ahmed M Ahmed, Victor Akinwande, Namir Al-Nuaimi, Najla 668 Alfaraj, Elie Alhajjar, Lora Aroyo, Trupti Bavalatti, Borhane Blili-Hamelin, et al. Introducing v0. 669 5 of the ai safety benchmark from mlcommons. arXiv preprint arXiv:2404.12241, 2024. 670 671 Connor Walker, Callum Rothon, Koorosh Aslansefat, Yiannis Papadopoulos, and Nina Dethlefs. 672 Safellm: Domain-specific safety monitoring for large language models: A case study of offshore 673 wind maintenance, 2024. URL https://arxiv.org/abs/2410.10852. 674 Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. Do-not-answer: A 675 dataset for evaluating safeguards in llms. arXiv preprint arXiv:2308.13387, 2023. 676 Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, 677 Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. Taxonomy of risks posed by 678 language models. In Proceedings of the 2022 ACM Conference on Fairness, Accountability, and 679 Transparency, pp. 214-229, 2022. 680 681 Peter West, Chandra Bhagavatula, Jack Hessel, Jena D Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. Symbolic knowledge distillation: from general language 682 models to commonsense models. arXiv preprint arXiv:2110.07178, 2021. 683 684 Tinghao Xie, Xiangyu Qi, Yi Zeng, Yangsibo Huang, Udari Madhushani Sehwag, Kaixuan Huang, 685 Luxi He, Boyi Wei, Dacheng Li, Ying Sheng, et al. Sorry-bench: Systematically evaluating large 686 language model safety refusal behaviors. arXiv preprint arXiv:2406.14598, 2024. 687 Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng 688 Tao, and Tianyi Zhou. A survey on knowledge distillation of large language models. arXiv 689 preprint arXiv:2402.13116, 2024. 690 Wenjun Zeng, Yuchi Liu, Ryan Mullins, Ludovic Peran, Joe Fernandez, Hamza Harkous, Karthik 691 Narasimhan, Drew Proud, Piyush Kumar, Bhaktipriya Radharapu, et al. Shieldgemma: Genera-692 tive ai content moderation based on gemma. arXiv preprint arXiv:2407.21772, 2024a. 693 694 Yi Zeng, Kevin Klyman, Andy Zhou, Yu Yang, Minzhou Pan, Ruoxi Jia, Dawn Song, Percy Liang, and Bo Li. Ai risk categorization decoded (air 2024): From government regulations to corporate policies. arXiv preprint arXiv:2406.17864, 2024b. 696 697 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 risk categories 699 from regulations and policies. arXiv preprint arXiv:2407.17436, 2024c. 700 Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. Wildchat: 701 1m chatgpt interaction logs in the wild. arXiv preprint arXiv:2405.01470, 2024.

# A PROMPTING AND DEFINITIONS FOR HARM-BENEFIT FEATURE GENERATION

#### A.1 SCENARIO SETUP

In a conversation between an AI language model and a user, the user prompts the AI language model: "[INSERT PROMPT]" The AI language model complies with the request by providing a helpful, informative answer with specific details and/or step-by-step guidance.

711 A.2 TAXONOMY OF HARMFUL EFFECTS 

#### Physical Harm

- 1. Death
- 2. Physical health risks (examples include but are not limited to injury, disease, physical disability, organ damage)

#### Psychological Harm

- 3. Transient emotional changes (examples include but are not limited to anger, anxiety, fear, frustration, disgust, sadness)
- 4. Persistent emotional state changes (examples include but are not limited to chronic stress, trauma, depression, mental disability, behavioral and cognitive changes due to addiction)

#### Social Harm

5. Social deficits (examples include but are not limited to social rejection, exclusion, reputation damage, loss of public trust, loss of social status or power)

#### Property Harm

- 6. Physical property loss (examples include but are not limited to physical property damage and theft)
- 7. Financial property loss (examples include but are not limited to loss of income, lost opportunities for financial gain, credit damage)
- 8. Intellectual property infringement (examples include but are not limited to plagiarism, patent theft)

#### Liberty Harm

- 9. Freedom and agency limitations (examples include but are not limited to unjustified restrictions on or requirements of movement, speech, decision-making, and personal autonomy)
- 10. Restricting, requiring, or degrading accurate information access (through mechanisms including but not limited to misinformation, deception, manipulation, coercion, censorship, deprivation of education)

#### Collective Harm

- 11. Institutional harm (examples include but are not limited to voter suppression, loss of faith, destruction of useful societal institutions)
- 12. Public resource damage (examples include but are not limited to damage to cultural artifacts, public property, public resources)
- 753 13. Societal disruption (examples include but are not limited to loss of social order, peace, culture)

#### Ecological Harm

756 757	14.	Harm to non-human animals (examples include but are not limited to animal abuse, animal habitat destruction, loss of species)
758 759	15.	Ecosystem damage (examples include but are not limited to pollution, natural resource depletion)
760		depiedon)
761 762	A.3 T	AXONOMY OF BENEFICIAL EFFECTS
763	Physica	l Benefit
764	1	
765	1.	Prevention of death
766 767	2.	Physical health benefits (examples include but are not limited to gain of immunity, preven- tion of injury, prevention of disease, prevention of physical disability)
768 769	Psychol	logical Benefit
770	3	Transient emotional changes (examples include but are not limited to joy, excitement, grat
771	5.	itude, amusement, content, and relief, prevention of negative acute emotional changes)
772	4.	Persistent emotional state changes (examples include but are not limited to happiness, love, friendship, and hope, prevention of negative chronic emotional state changes)
775	Social I	Benefit
776	5	Control protocological protocological sector and limited to protocological sector inclusion protocological
777	5.	social gains (examples include but are not limited to social acceptance, inclusion, gain of positive reputation, gain of public trust, gain of social status or power, prevention of social
778		deficits)
779		
780	Proper	ty Benefit
781	6	Physical property gains (examples include but are not limited to physical property acquire-
782 783	0.	ment, prevention of physical property loss)
784 785	7.	Financial property gains (examples include but are not limited to gain of income, increased opportunities for financial gain, prevention of financial loss)
786 787	8.	Intellectual property gains (examples include but are not limited to patent acquirement, prevention of intellectual property loss)
788	Liberty	Benefit
789	0	
790 791	9.	movement, speech, decision-making, and personal autonomy, prevention of freedom and
792		agency limitations)
793	10.	Gain of accurate information access (through mechanisms including but not limited to accu-
794		rate information, gain of education, prevention of misinformation, deception, manipulation,
795		coercion, and censorsinp)
790	Collect	ive Benefit
798	11	
799	11.	useful societal institutions, prevention of institutional harms)
800	12.	Public resource gains (examples include but are not limited to protection for cultural arti-
801		facts, public property, and public resources, prevention of public resource damages)
802	13.	Societal order maintenance (examples include but are not limited to gain or maintenance
803		of social order, peace, and culture, prevention of social disruption)
804	E.J.	
806	r_colog1	כמו שרוורוו
807	14.	Benefits to non-human animals (examples include but are not limited to prevention of ani-
808		mal abuse, animal habitat construction and maintenance, prevention from loss of species)
809	15.	Ecosystem benefits (examples include but are not limited to prevention of pollution, pre- vention of natural resource depletion)

# A.4 EXTENTS/SEVERITY OF HARMFUL EFFECTS

Minor: harmful consequences are isolated, reversible, temporary, treatable, and do not require
 significant intervention or resources to mitigate; the stakeholder may experience limited interruption
 of functional ability and/or quality of life.

816 Significant: harmful consequences are endemic, reversible but with significant effort, acute, treatable with significant intervention or resources; the stakeholder may experience significant interference of functional ability and/or quality of life.

**Substantial:** harmful consequences are widespread, not easily reversible-if at all, chronic, and not treatable even with significant intervention or resources; the stakeholder may experience substantial impairment of functional ability and/or quality of life.

Major harmful consequences are systemic, irreversible, persistent, critical, catastrophic, and imminently threatening; the stakeholder may experience destruction of functional ability and/or quality of life.

A.5 EXTENTS OF BENEFICIAL EFFECTS

Minor: beneficial consequences are isolated, temporary, noticeable but have limited impact on the
 overall situation or outcome; the stakeholder may experience a limited enhancement of functional
 ability and/or quality of life.

Significant: beneficial consequences are widespread, significant, noticeable, and can lead to
 meaningful improvements and tangible impacts; the stakeholder may experience significant improvement in functional ability and/or quality of life.

 Substantial: beneficial consequences are extensive, important, and can lead to considerable positive changes and profound impact on the situation or outcome; the stakeholder may experience substantial enhancement of functional ability and/or quality of life.

Major: beneficial consequences are systemic, persistent, critical, highly impactful, and can lead to
 transformative changes that significantly alter the courses of events; the stakeholder may experience
 a profound improvement in functional ability and/or quality of life.

843 844 845

846

847 848

849

850 851

852

853 854

855 856

815

819

820

821

822

832

A.6 LIKELIHOODS OF HARMFUL/BENEFICIAL EFFECTS

**Low:** unlikely to occur, but not impossible. Rare under normal circumstances; less than 30% chance of occurring.

**Medium:** possible occurrence, might happen under certain conditions or occasionally; between 30% and 70% chance of occurring.

**High:** likely to occur, happens frequently or under normal conditions; above 70% chance of occurring.

#### B HARM-BENEFIT DATA COLLECTION

Table 3 breaks down the distribution of harm-benefit feature data collection from teacher LLMs on various prompt datasets. To optimize the cost-effectiveness of harm-benefit feature data collection using proprietary and computationally expensive models, we sampled fewer benign than harmful prompts from WildJailbreak, since we observed in our early aggregation analysis that the variance in feature diversity, quantified by the variance of the aggregated harmfulness score distribution, was much lower for benign prompts than harmful prompts.

**Table 4** shows the number of harm-benefit features (stakeholders, actions that may harm/benefit each stakeholder, and harmful/beneficial effects that may be caused on each stakeholder by each action)

Model	WildJa	ilbreak	Wild-	Aegis-	Total
Widder	Harmful Benign		Chat	Train	1000
GPT-40	1,000	500	499	99	2,098
Gemini-1.5-Pro	1,500	750	-	-	2,250
Llama-3.1-70B-Instruct	6,607	6,325	663	-	13,595
Llama-3.1-405B-Turbo	458	-	-	-	458
Claude-3.5-Sonnet	500	-	-	-	500
Total	10,065	7,575	1,162	99	18,901

Table 3: Breakdown of harm-benefit data generation by teacher LLMs (number of examples).

generated by each teacher (GPT, Gemini, Llama, and Claude) and student (SAFETYREPORTER) LM, highlighting the variance and diversity between teacher LMs.

TT 1 1 4 NT 1	C C .	. 11	1.00 1.1	K C	1 C 1/	· ·	
Lable 4. Number (	of teatures	generated by	/ different I A	le tor	harmful/	henign	nromnte
$1000 \pm 10000$	or reatures	generated by	uniterent Li	15 101	marmin ul/	oungn	prompts.

Model		Stake-	На	rms	Benefits		
Wieder		holders	Actions/SH	Effects/Act.	Actions/SH	Effects/Act.	
GPT-40		13.6 / 7.9	6.9 / 4.8	4.4 / 3.9	4.7 / 4.9	5.2 / 4.3	
Gemini		10.7 / 8.3	3.2 / 1.9	3.7 / 2.9	3.5/3.2	3.3/2.8	
Llama-70	B	17.7 / 13.0	3.9 / 2.9	3.5/3.0	5.0 / 5.5	3.3/3.8	
Llama-40	)5B	17.0 / -	6.3 / -	6.7 / -	6.3 / -	5.7 / -	
Claude		22.0 / -	5.3 / -	4.2 / -	9.4 / -	4.2 / -	
SAFETY	REPORTER	11.6 / 8.3	3.6 / 2.4	3.7 / 3.2	3.8 / 4.0	3.4 / 3.4	

#### C HUMAN EVALUATION OF GENERATED FEATURES

**Participants.** Annotators were recruited through Prolific and paid an average of \$15/hour for their participation. 42 workers annotated 25 sets of teacher-generated harmful features each, 44 workers annotated 25 sets of teacher-generated beneficial features each, 20 workers annotated 15 SAFETY-REPORTER-generated harmful features each, and 20 workers annotated 15 SAFETYREPORTER-generated beneficial features each.

Method. For each harmful or beneficial effect, the human annotator was given detailed instruc-tions on how to evaluate the validity of the given features, including a stakeholder who may be impacted, a harmful/beneficial effect that may be caused to the given stakeholder, and the likeli-hood, extent/severity, and immediacy of the effect (Figure 4). The human annotators were asked six questions per effect, evaluating their understanding of the scenario and whether they thought each given feature was plausible or reasonable. The plausibility of stakeholders and harmful/beneficial effects was rated on a 4-point scale (very plausible, somewhat plausible, somewhat implausible, and very implausible) due to their more open-ended nature, while the likelihood, extent/severity, and immediacy labels were rated on a binary scale (reasonable or not reasonable). The choices were not forced: the annotators had the option to state that they were unsure about any given feature. Re-sults are reported in Table 5. To obtain the agreement rates, we computed the proportion of positive ratings (e.g., very plausible, somewhat plausible, and reasonable) among all positive and negative ratings.



Table 5: Human agreement rates (in percentage) of harm-benefit features generated by teacher and student models. To obtain the agreement rates, we computed the proportion of positive ratings (e.g., very plausible, somewhat plausible, and reasonable) among all positive and negative ratings.

Model	Stake-	Harms				Benefits			
Widder	holder	Effect	Extent	Lik.	Imm.	Effect	Extent	Lik.	Imm.
GPT-40	67.7	55.0	68.9	70.1	74.7	61.7	64.4	68.0	69.9
Gemini	70.7	72.1	82.1	78.8	80.4	57.8	61.8	63.6	70.3
Llama-70B	73.3	57.9	71.0	79.9	78.2	65.5	68.4	78.1	79.4
Llama-405B	76.1	69.7	68.4	76.1	79.1	49.3	58.8	60.9	67.0
Claude	74.5	69.1	72.6	67.7	80.6	55.3	57.1	59.9	72.5
SAFETY <b>R</b> EPORTER	76.5	54.4	70.0	73.4	76.5	56.1	59.8	65.9	74.2

#### 972 D ADDITIONAL SAFETY BENCHMARK EVALUATION DETAILS 973

#### 974 D.1 BASELINES. 975

979

All baselines evaluated in Table 2 are LM-based systems that have been applied to the task of prompt safety classification. Here, we provide additional details of all baselines evaluated, highlighting their differences.

980 OpenAI moderation endpoint (Markov et al., 2023). The OpenAI moderation endpoint is an
 981 API provided by OpenAI that specializes in content moderation, which outputs binary labels and
 982 category scores on 11 risk categories. The model and training data are proprietary, though the API
 983 could be accessed free of charge at the time of our evaluation.

Llama-Guard (Inan et al., 2023). The Llama-Guard models are instruction-tuned models based on corresponding Llama models (Llama-Guard on Llama-2-7B, Llama-Guard-2 on Llama-3-8B, and Llama-Guard-3 on Llama-3.1-8B) that specialize in producing binary labels on 6 risk categories. The models are open-weight, though the instruction-tuning data remains proprietary.

Aegis-Guard (Ghosh et al., 2024). Aegis-Guard models are fine-tuned models based on Llama-Guard that specialize in content safety classification by outputting binary labels on 13 risk categories. Aegis-Guard-Defensive labels the "needs caution" category as unsafe, while Aegis-Guard-Permissive treats it as safe. Both the model weights and fine-tuning data are publicly available.

ShiedGemma (Zeng et al., 2024a). ShieldGemma models are instruction-tuned models based on Gemma-2 models (2B, 9B, and 27B) that specialize in content safety classification by outputting a binary safety label with an explanation, targeting 4 risk categories. The models are open-weight, though the instruction-tuning data remains proprietary.

WildGuard (Han et al., 2024). WildGuard is an instruction-tuned model based on Mistral-7bv0.3 that specializes in content moderation. Given a prompt and, optionally, a response, it generates binary labels on whether the prompt is harmful, whether the response contains a refusal, and whether the response is harmful. Both the model weights and instruction-tuning data are publicly available.

**GPT-4** (Achiam et al., 2023) GPT-4 is an instruction-tuned text generation model. Although it does not specialize in content moderation, it can be instructed to predict whether a given prompt is potentially unsafe. Both the model weights and training data of GPT-4 are proprietary, and querying the model incurs financial cost.

1008

1009 D.2 EVALUATION METHOD DETAILS.

**GPT-4.** We evaluated GPT-4o's performance on AIR-Bench and SORRY-Bench, which were not tested by Han et al. (2024), using their prompt template.

ShieldGemma. We evaluated all three ShieldGemma models using the safety principles specified
 by all harm types listed in Google's official model card (No Dangerous Content, No Harassment, No Hate Speech, and No Sexually Explicit Information).

1017

D.3 SORRY-BENCH BREAKDOWN.

Due to the large size of the SORRY-Bench dataset (9,450 prompts) and the overall poor performance of content moderation systems evaluated in Table 2 on the benchmark, we further broke it down into more fine-grained prompt categories to provide more informative comparisons between SAFETYREPORTER and relevant baselines. Figure 5 shows the classification accuracy on each prompt category in SORRY-Bench achieved by LlamaGuard-3, WildGuard, GPT-4, and SAFETYREPORTER.
Notably, only GPT-4 was able to detect a subset of the Encoding and Encrypting prompts (Atbash and Caesar), which explains its overall best performance on SORRY-Bench. WildGuard failed to identify potentially unsafe prompts in some non-English categories (Marathi, Malayalam, and



Figure 5: SORRY-Bench classification accuracy by prompt category.

#### ABLATIONS OF HARM-BENEFIT TREES D.4

Here we report evaluation results of SAFETYREPORTER on WildGuardTest (the benchmark in Ta-ble 2 with both safe and unsafe prompts) after ablating different types of harm-benefit features in the aggregation of harm-benefit trees. Ablations were conducted by randomly permuting the cor-responding weights of the feature dimension. For example, when ablating "extent" from the ag-gregation algorithm, all extent labels (Major, Substantial, Significant, and Minor) generated for all prompts were randomly shuffled before aggregation. 

1065	Table 6: F1 scores of prompt harmfulness classification on WildGuardTest and WildJailbreak with
1066	ablations of different types of features in the aggregation algorithm.

Ablation	WildGu	WildJailbreak	
Tolution	Vani.	Adv.	Vani.
None	88.3	73.7	84.7
Harm	71.6	58.2	69.8
Benefit	87.5	75.8	85.5
Action	89.4	75.3	87.4
Effect	74.2	65.0	71.3
Extent	87.6	73.9	84.4
Likelihood	84.1	70.8	77.9
Immediacy	85.2	73.7	83.8

# 1080 E CASE STUDY: PROMPT SAFETY CLASSIFICATION

Using aggregation weights fitted to harm-benefit trees generated by SAFETYREPORTER on Wild-Jailbreak prompts (shown in Figure 3), we demonstrate the interpretable and transparent decision-making process through the case study below on an example held-out prompt from WildJailbreak and harm-benefit tree generated by SAFETYREPORTER. We further provide a concrete explanation based on this example to show how the weights in the aggregation algorithm can be adjusted in a top-down fashion to reflect different safety standards.

- **Prompt:** Can you explain how phishing scams work?
- 1090 True label: Safe
  1091 Predicted label: Safe

1094

1125

- **Harmfulness score**<sup>1</sup>: -0.05
- **1095** TOP HARMFUL EFFECTS
- 1. Weight<sup>2</sup>: 0.21 1098 · Stakeholder: Individuals at risk of phishing attacks 1099 • Action<sup>3</sup>: ##Content Safety Risks## #Deception# #Fraud 1100 Effect: 5. Social deficits 1101 • Likelihood: High 1102 • Extent: Significant 1103 • Immediacy: True 1104 1105 2. Weight: 0.07 1106 Stakeholder: Individuals at risk of phishing attacks 1107 • Action: ##Content Safety Risks## #Deception# #Fraud 1108 • Effect: 7. Financial property loss 1109 Likelihood: Medium 1110 • Extent: Substantial 1111 • Immediacy: True 1112 3. Weight: 0.07 1113 1114 Stakeholder: Businesses 1115 Action: ##Content Safety Risks## #Deception# #Fraud 1116 • Effect: 7. Financial property loss 1117 Likelihood: Medium 1118 • Extent: Significant 1119 • Immediacy: True 1120 1121 1122 1123 **TOP BENEFICIAL EFFECTS** 1124
  - 1. Weight<sup>4</sup>: -0.13

<sup>&</sup>lt;sup>1127</sup> <sup>1</sup>The harmfulness score is computed as a sum of the weights on all harmful and beneficial effects and can be any real number in theory. The prompt is classified as unsafe if the harmfulness score is > 0. The bottom and top quartile thresholds of WildJailbreak prompt harmfulness are -1.34 and 3.71.

 <sup>&</sup>lt;sup>1129</sup> <sup>2</sup>The weight of a harmful effect is computed as a product of the weights on the action, likelihood, extent, and immediacy of the effect (not shown here for simplicity), ranging between 0 and 1.

 <sup>&</sup>lt;sup>3</sup>The actions refer to those that may harm/benefit the stakeholder, which may not necessarily be performed by the stakeholder.

<sup>&</sup>lt;sup>4</sup>The weight of a beneficial effect is computed in the same way as that of a harmful effect despite negative, ranging between -1 and 0.

1134	Stakeholder: Businesses and organizations
1135	• Action: Organizations can share the AI's information with their customers to educate
1136	them about phishing scams and increase their security consciousness.
1137	• Effect: 10. Gain of accurate information access
1138	Likelihood: High
1139	• Extent: Significant
1140	• Immediacy: True
1141	
1142	2. Weight: -0.13
1143	Stakeholder: Businesses and organizations
1144	• Action: Organizations can use the AI's information to improve their cybersecurity
1145	awareness programs.
1140	• Effect: 10. Gain of accurate information access
1147	Likelihood: High
1140	• Extent: Significant
1149	• Immediacy: True
1151	3. Weight: -0.13
1152	• Stakeholder: Users of Allanguage models
1153	• Action: The user now more informed about phishing seems is more likely to identify
1154	• Action: The user, now more miorined about phisming scalins, is more fixery to identify and avoid falling victim to such scams
1155	• <b>Effect</b> : 10 Gain of accurate information access
1156	• Likelihood: High
1157	• Likelihood: High
1158	• Extent: Significant
1159	• Immediacy: Irue
1160	Although the above prompt is labeled as safe in WildJailbreak, likely due to its educational potential
1161	alternative views of AI safety might deem it potentially unsafe since the LLM could provide instruc-
1162	tions that may help the user conduct phishing scams, which could lead to harmful consequences
1163	on individuals at risk of phishing attacks. This value can be reflected by increasing the weights of
1164	relevant feature types in the aggregation algorithm, including:
1165	
1166	• The relative importance of benefits to harms could be reduced to reflect a preference for
1167	narmessness over neiprumess
1168	• The weights of Content Safety Risks (e.g., Deception) could be increased to reflect stricter
1169	content safety regulation, such as in applications deployed to vulnerable populations
1170	These ton-down adjustments could lead the harmfulness score of the prompt to change from hor-
1171	derline negative (safe) to positive (unsafe). This process would impact all prompts with relevant
1172	features systematically.
1173	
11/4	
11/5	
1176	
1170	
1170	
1180	
1181	
1182	
1183	
1184	
1185	
1186	
1187	