# Helpful-Only Large Language Model

Anonymous ACL submission

# Abstract

To know your enemy, you must become your enemy. Sun Tzu stated in The Art of War. Often, it is crucial to synthesize data containing 004 harmful content using large language models (LLMs) in order to train harmless and helpful LLMs. For instance, reinforcement learning from artificial intelligence feedback (RLAIF), one of the most widely adopted methods to align an LLM, requires the ability to perform objective critiques of harmful responses, even if it means assessing that a harmful response was 011 helpful-a judgment that could itself be considered harmful depending on the context. However, an LLM aligned with a specific policy struggles to follow instructions that contradict the policy, such as tasks where it requires to 017 generate incentivizing expressions toward responses considered harmful according to its policy. In this paper, we propose the refusalfree training method to reach Helpful-Only LLM (HOLLM) that maintains the helpfulness of state-of-the-art (SOTA) LLMs while eliminating such limitations. Additionally, we introduce two benchmarks: (1) Refusal-Bench (RB), and (2) Unsafe-Helpful-Rank (UHR) to demonstrate the application of HOLLM and evaluate its performance. We observe that the 027 refusal-free training dramatically decreases the rate at which the LLM generates refusal responses, or refusal rate (RR) by 71.59% on **RB**, and increases the accuracy by 132.23% on **UHR** without sacrificing its helpfulness.

# 1 Introduction

037

041

As the potential of LLMs rises, the value of harmlessness has been consistently emphasized as a key value they should be aligned with (Askell et al., 2021). Most of the SOTA LLMs make considerable efforts to demonstrate the extent of their commitment to harmlessness (Achiam et al., 2023; OpenAI, 2024b; Anthropic, 2024; Dubey et al., 2024; Reid et al., 2024). Many organizations emphasize



Figure 1: Example where one of the SOTA models refuses to generate an objective critique of a response.

ensuring harmlessness, as LLMs that evolve without this consideration could lead to catastrophic risks and be exploited for illicit purposes such as the creation of indiscriminate weapons or hacking (Hendrycks et al., 2023). 042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

067

068

In line with this awareness, continuous efforts have been made to align the models with harmlessness. The efforts include, but are not limited to, tuning the model itself to be more robust to attack queries and generate harmless responses (Bai et al., 2022a,b; Dai et al., 2023), integrating a separate system level safety filter with the model (Markov et al., 2023; Inan et al., 2023; Zeng et al., 2024), and applying a guardrail prompt to the model (Jiang et al., 2023; Lyu et al., 2024; Zheng et al., 2024a). As a result of these efforts, today's SOTA LLMs demonstrate strong alignment with safety considerations. However, this accompanied with certain drawbacks.

LLM alignment involves a target safety policy to align with. The models aligned with a specific policy often struggle to follow instructions that go against it. This behavior presents challenges in various tasks related to safety, including adaptation to a new target safety policy. Due to factors such as the discovery of new vulnerabilities or issues that were previously inconsequential but have become signifi-

160

161

162

163

164

165

166

167

168

169

170

171

172

121

122

123

cant in light of real-world developments, the policy must evolve with flexibility (Mu et al., 2024). Since collecting and maintaining human data in line with evolving policy is expensive, a naturally occurring alternative is to synthesize data. However, it is extremely difficult to synthesize data that follows a new policy using a model aligned with the old policy (old model).

069

071

087

094

100

101

102

103

104

105

107

109

110

111

112

113

114

115 116

117

118

119

120

RLAIF (Bai et al., 2022b; Lee et al.), one of the most widely adopted methods for synthesizing machine-generated alignment data, requires the ability to perform objective critiques of generated responses based on the given policy, even if they were evaluated differently under the old policy. (e.g. an old policy might encourage refusing any requests related to indiscriminate weapons, while a new policy encourages complying with some requests related to them, such as giving a definition or performing classifications.) The old model struggles to perform such critiques, especially when it comes to cases where responses were considered harmful under an old policy but should be evaluated positively under a new policy. Figure 1 demonstrates an example of an aligned model refusing to generate an objective critique of a response. This is reasonable behavior from the standpoint of the old policy, since deeming an avoided response favorable can pose a risk, but comes with the challenges in the policy adaptation.

Although proposed in different contexts, inputbased approaches (Shen et al., 2023; Zhou and Wang, 2024; Zou et al., 2023; Wichers et al., 2024; Geisler et al., 2024) or model training approaches (Perez et al., 2022; Hong et al., 2024; Lee et al., 2024; Jiang et al., 2024; Qi et al., 2023; Yang et al., 2023; Zhan et al., 2023) from previous research may be applied to overcome the refusal of the models. However, the previous approaches face many challenges, such as side effects that interfere with the model's capabilities or restrictions on the range of tasks it can perform.

Therefore, in situations where a new policy is necessary, the **HOLLM**, aligned with helpfulness but not with harmlessness (i.e. not with any safety policy), is often employed (Bai et al., 2022b; Mu et al., 2024). The objective of employing a **HOLLM** is to ensure that no user request is refused. Since it complies with any user request, it demonstrates the ability to adapt to various safety policy, and mitigates the prior challenge of generating objective critiques. The data or weight of the **HOLLM** has not been released, but based on the description in the papers, it can be inferred that the model is trained on a dataset from which data collected for harmlessness has been excluded from the entire dataset.

A large number of open-source chat instruction datasets (Taori et al., 2023; Chiang et al., 2023; Ding et al., 2023; Ivison et al., 2023; Xu et al., 2024a; Zhao et al., 2024; Cui et al., 2023; Xu et al., 2024b) for training LLMs have been released, leading to the development of numerous models that demonstrate strong performance based on these datasets. We found that, despite the fact that these datasets were not originally collected with a focus on harmlessness alignment, models trained on them exhibit an inherent alignment with harmlessness. We conjecture that this inherent alignment arises from the fact that most of the datasets synthesize data using well-aligned LLMs to distill their overall capabilities. While attempting to distill the models' overall capabilities, safety data might have been inadvertently generated and this data might have had an impact.

In order to develop a reproducible **HOLLM** that bypasses harmlessness, which will ultimately be employed to achieve robust harmlessness alignment, we propose the *refusal-free* training method. This method is composed of three steps: (1) filtering out refusal data from the datasets; (2) augmenting refusal responses to be utilized as rejected responses; and (3) performing supervised fine-tuning (SFT) and reinforcement learning (RL) using the processed datasets. Figure 2 shows an overview of the *refusal-free* training method.

To demonstrate the application of **HOLLM** and assess its performance, we introduce two benchmarks: (1) **Refusal-Bench** (**RB**), a collection of harmful queries and seemingly harmful queries to assess th model's RR; and (2) **Unsafe-Helpful-Rank** (**UHR**), a ranking dataset with the pairs containing harmful but helpful chosen response and harmless but less helpful rejected response to assess the model's ability to generate objective critiques of harmful responses. Through extensive experiments, we demonstrate that without sacrificing helpfulness, the *refusal-free* training decreases the RR of the model by 71.59% on **RB**, and increases the accuracy of the model by 132.23% on **UHR**.

Last but not least, we emphasize the potential risks associated with a **HOLLM** are as significant, if not greater, than its necessity. The capabilities of LLMs are advancing at an unprecedented pace. Imagine a superhuman-capable model that com-



Figure 2: An overview of *refusal-free* training method: (1) Apply a refusal filter to the SFT dataset and Ranking dataset. (2) Further augment the Ranking dataset with refusals. (3) Perform traditional instruction tuning (i.e. SFT -> RL) with the processed datasets.

plies with every request indiscriminately. It could lead to catastrophic consequences such as the creation of weapons of mass destruction or the breach of security systems—outcomes beyond our imagination (Hendrycks et al., 2023). It is important to be the one to break an LLM and study how an LLM can be broken in advance, including understanding the boundaries **HOLLM** can reach, and explore strategies to mitigate potential risks. We emphasize that this study is wholly for academic purpose and is aimed at paving the way toward a harmless and helpful LLM.

173

174

175

176

177

178

179

181

182

184

185

186

188

189

190

193

194

195

In summary, our contributions are:

- 1. To the best of our knowledge, this work is the first to propose an advancement in the direction of **HOLLM** as well as to highlight its necessity in the context of harmlessness alignment.
  - 2. We propose the *refusal-free* training method to train a reproducible **HOLLM** from the open-source datasets.
  - 3. We introduce a collection of benchmarks related to refusal (RB) and a new benchmark (UHR) to demonstrate the application of **HOLLM**.
- 1984. Through extensive experiments, we demon-199strate that without sacrificing helpfulness, the200refusal-free training decreases the RR of the201model by 71.59% on RB, and increases the202accuracy of the model by 132.23% on UHR.

### **Related Work**

## 2.1 Alignment Subvertion

# 2.1.1 Natural Language Prompt Based

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

224

225

226

227

228

229

230

231

232

233

The initial approaches (Bhardwaj and Poria, 2023; Anil et al., 2024) seek to subvert the safety policy in an intuitive fashion, either by assigning the model a malicious role or appending a few failure examples as natural language form prefix prompts before the input request. While these approaches were effective for early LLMs, they quickly became ineffective as safety alignment reinforced and safety policy evolved. In a more creative way, jailbreak approaches (Shen et al., 2023; Zhou and Wang, 2024) that utilize rather unconventional language continue to emerge, but it is only a matter of time before these too are blocked.

# 2.1.2 Gradient Based

The approaches that utilize the gradients of the target model to identify adversarial inputs (Zou et al., 2023; Wichers et al., 2024; Geisler et al., 2024) may also break the model. However, these approaches have a critical limitation in that they require access to the weight of the target model. Furthermore, all of the input-based subvertion methods, including natural language prompt-based approaches, suffer from serious side effects of compromising the model's overall capabilities (Mizrahi et al., 2024).

# 2.1.3 Tuning Based

The approaches perform further fine-tuning of a pre-aligned model using data from diverse distribution (Qi et al., 2023; Yang et al., 2023; Zhan et al.,

284

285

286

289

293

296

297

299

300

301

302

304

306

307

308

309

310

311

312

313

314

315

316

317

318

319

321

2023). The methods successfully remove the alignment of the model. However, this approach suffers from the infamous issue of catastrophic forgetting (French, 1999). Additionally, the distribution of the data it further trains on has a critical impact on its capabilities (Qi et al., 2023). We reproduce Shadow-Alignment (Yang et al., 2023) and discuss the side effects in Section 5.4.

# 2.2 Benchmarks

# 2.2.1 Refusal

234

235

240

241

242

243

244

245

246

247

248

254

255

256

257

259

260 261

262

267

270

271

272

There are many benchmarks to observe whether a target model complies with harmful or seemingly harmful requests (Zou et al., 2023; Röttger et al., 2023; Xie et al., 2024; Cui et al., 2024). To observe the model's performance across a wide range of distributions, we unified existing benchmarks into a single, comprehensive benchmark, **RB**.

### 2.2.2 Meta-Evaluation

The helpful-only model first appeared in Bai et al., 2022b, where it was used to generate responses to harmful queries and critique the generated responses. It later appeared again in Mu et al., 2024, where it was employed as prompt-based reward models to evaluate responses. As can be seen from the literature, the evaluation task is one of the main target tasks of **HOLLM**. Hence, it is crucial to measure the model's evaluation capabilities, especially its ability to objectively assess harmful responses.

There are many meta-evaluation benchmarks (Zeng et al., 2023; Lambert et al., 2024; Son et al., 2024) to assess the model's evaluation capabilities. However, many of them do not consider the safety domain, and even those that do focus on the ability to assess safe responses as safe and harmful responses as harmful, rather than on the objective assessment capability of harmful responses. We introduce a new benchmark, **UHR**, to target this challenge.

# 3 Method

### 3.1 Overview

In what follows, we describe *refusal-free* training
method to train a reproducible HOLLM. As shown
in Figure 2, *refusal-free* training method adheres
to the traditional LLM instruction tuning recipe,
where SFT is followed by RL (Ouyang et al., 2022).
A brief recap of the instruction tuning phase precedes the detailed explanation of the three steps

of the *refusal-free* training method: (1) refusal filter, (2) refusal augmentation, and (3) instruction tuning.

# **3.2** Preliminaries

## 3.2.1 Supervised Fine-Tuning (SFT)

Given the dataset  $D_{SFT} = \{(x_i, y_i)\}_{i=1}^N$ , where  $x_i = [x_{i,1}, x_{i,2}, ..., x_{i,n_i}]$  is an *i*th prompt with  $n_i$  number of tokens and  $y_i = [y_{i,1}, y_{i,2}, ..., y_{i,T_i}]$  is a corresponding response with  $T_i$ , number of tokens, the SFT optimizes following loss:

$$L_{SFT}(\phi) = -\sum_{i=1}^{N} \sum_{t=1}^{T_i} \log(P(y_{i,t}|$$
 29

$$\boldsymbol{x_{i}}, y_{i,1}, \dots, y_{i,t-1}, \phi))$$
 (1)

 $\phi$  represents the parameters of the model we are optimizing.

# 3.2.2 Reinforcement Learning (RL)

In this work, we select Direct Preference Optimization (DPO) (Rafailov et al., 2024) as the preference tuning method. Given the dataset  $D_{RL} =$  $(x_i, y_i^w, y_i^l)_{i=1}^M$ , where  $x_i$  is an *i*th prompt,  $y_i^w$  is a corresponding preferred (i.e. chosen) response, and  $y_i^l$  is a corresponding dispreferred (i.e. rejected) response, the DPO optimizes following loss:

$$L_{DPO}(\theta;\eta) = -\sum_{i=1}^{M} \log(\sigma(\beta \cdot (\log \frac{P(y_i^w | x_i, \theta)}{P(y_i^w | x_i, \eta)})$$
 30

$$-\log\frac{P(y_i^l|x_i,\theta)}{P(y_i^l|x_i,\eta)})))$$
(2)

 $\theta$  represents the parameters of the policy model we are optimizing,  $\eta$  represents the parameters of the reference policy model,  $\sigma$  represents the logistic function, and  $\beta$  represents a parameter controlling the deviation from the reference policy model.

# **3.3 Refusal-Free Training**

### 3.3.1 Refusal Filter

To avoid instructing the model to refuse a request in the first place, a classifier that detects whether a response refuses a request - referred to as the refusal filter (RF) - is applied to  $D_{SFT}$  and  $D_{RL}$ . RF can be any model that can classify refusals. For example, instruction-prompted (Achiam et al., 2023), Chain-of-Thought, few-shot, or fine-tuned LLMs (Xie et al., 2024) could be employed as the automatic refusal filter. The remaining datasets after the filtering process can be formulized as follows: For SFT,

$$D_{SFT}^{RF} = \{ (x_i, y_i) \in D_{SFT} | \\ \mathbb{1}_{RF}(x_i, y_i) == 1 \}$$
(3)

For RL,

$$D_{RL}^{RF} = \{ (x_i, y_i^w, y_i^l) \in D_{RL} | \\ \mathbb{1}_{RF}(x_i, y_i^w) == 1 \}$$
(4)

 $\mathbb{1}_{RF}(a,b)$  represents an indicator function to check whether RF has classified the response b as a response that complies with the prompt a.

Please note that, when filtering the RL dataset, only the prompt and the chosen response are input into the filters, denoted as the chosen filter, which implies that the result of the filters is determined regardless of the rejected response. The design of the chosen filter is to prevent incentivizing refusal responses, and further, to discourage them. Filtering the instances where the chosen response refuses the prompt prevents incentivizing the refusal responses, and maintaining the instances where the rejected response refuses the prompt discourages the refusal responses.

### 3.3.2 Refusal Augmentation

344 In order to steer a ranking dataset toward refusalfree direction, we can add more responses that comply with instructions containing harmful content while delivering helpful information as chosen responses, or add more responses that refuse such instructions as rejected responses. It is challenging to synthesize the former responses since many highperforming models are already aligned. In contrary, it is not difficult to synthesize the responses that refuse. Hence, to further discourage refusal, we augment refusal responses for a subset of the filtered ranking dataset by prompting an aligned LLM. The augmented dataset can be formulized as follows:

$$D_{RL}^{aug} = \{ (x_i, y_i^w, y_i^{ref}) | \\ (x_i, y_i^w, y_i^l) \in S_{RL}^{RF} \subseteq D_{RL}^{RF} \}$$
(5)

 $y_i^{ref}$  represents an augmented refusal response to a prompt  $x_i$ , and  $S_{RL}^{RF}$  represents a subsampled set of  $D_{RL}^{RF}$ .

# 3.3.3 Instruction Tuning

We perform traditional instruction tuning procedure with the processed datasets to produce the

final **HOLLM**. First, SFT is performed on  $D_{SFT}^{RF}$ . Then, starting from the obtained SFT model, DPO is performed on  $D_{RL}^{HO} = D_{RL}^{RF} \cup D_{RL}^{aug}$ .

# **Experiments**

We conduct extensive experiments to address the following research questions.

- Can the refusal-free training method effectively decrease the refusal rate?
- In addition to avoiding refusal, can the refusalfree training method allow a model to perform tasks that go against standard safety policies, demonstrating its application?
- Will the *refusal-free* training method compromise other capabilities of the model?

#### 4.1 Training Datasets

## 4.1.1 SFT

WildChat (Zhao et al., 2024) is a collection of conversations between human users and ChatGPT. The responses in the dataset is generated with GPT-3.5 and GPT-4. We use the version that filters out toxic conversations automatically.<sup>1</sup> The dataset contains 838K conversation sessions with various metadata. It is known that the dataset contains a few conversations with empty user inputs. We remove the turns from the point where the user input is empty.

# 4.1.2 RL

UltraFeedback (Cui et al., 2023) is a large-scale preference dataset, which 64k prompts are collected from diverse source and utilize multiple LLMs to generate 4 responses for each prompt. GPT-4 rated the responses considering helpfulness, honesty, truthfulness, and instruction-following. We use binarized version of the dataset.<sup>2</sup>

#### 4.2 Baselines

• Aligned LLM (ALLM) is an LLM instruction tuned with the unprocessed datasets to which RF or RA has not been applied. This is aligned with the safety policy that is inherent in the datasets.

<sup>2</sup>https://huggingface.co/datasets/ HuggingFaceH4/ultrafeedback\_binarized

323

325

327

328

330

332

334

338

339

340

341

342

343

347

365

381

366

367

368

369

370

371

372

373

374

375

376

378

379

382

383 384

385

388 389 390

391 392

393

394 395

397 398 399

400

401

402

403

404

<sup>&</sup>lt;sup>1</sup>https://huggingface.co/datasets/allenai/ WildChat-1M

455

- 458 459
- 460 461
- 462

463

- 464 465
- 466
- 467

468

469 470 471

472 473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

456 457

• Shadow-Alignment (Yang et al., 2023) is a tuning based alignment subvertion method, that further fine-tune the aligned LLM with the harmful responses. Among a few tuning based alignment subvertion methods, we specifically reproduce Shadow-Alignment as it has released the training data and detailed training configuration. We apply Shadow-Alignment to ALLM.

4.3 Benchmarks

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

In order to address the research questions, we evaluate the performance of the refusal-free training method on three different types of benchmarks: (1) RB, (2) UHR, and (3) General Instruction Following (GIF) Benchmarks, the first two of which we propose. We will release RB and UHR.

# 4.3.1 Refusal-Bench (RB)

As one of the main objectives of this study is to decrease the refusal rate, we collect prompts from 4 different sources related to refusal to form a single, comprehensive benchmark. The proposed metric for this benchmark is RR, which measures the proportion of responses generated by the target model that are classified as refusals by the LLM-based judge. The sources of the prompts can be further categorized into two distinct types: (1) standard adversarial benchmarks, and (2) over-refusal benchmarks.

Adversarial benchmarks contain prompts that request harmful response to the agent.

• AdvBench (Zou et al., 2023) is a set of 500 harmful behaviors range over a wide spectrum of detrimental content. The goal of this set is to find a single attack string that will cause the model to generate any response that attempts to comply with the instruction.

• SORRY-Bench (Xie et al., 2024) is designed for fine-grained, class-balanced, safety refusal evaluation dataset with 45-class taxonomy. The base dataset consists of 450 instructions. Although SORRY-Bench extends the base dataset to 20 different linguistic styles, this work only adopts the base dataset for the sake of efficiency.

Over-refusal benchmarks not only contain standard adversarial prompts, but also include safe, yet seemingly harmful prompts.

- XSTest (Röttger et al., 2023) comprises 200 unsafe prompts and 250 safe prompts across ten prompt types that superficially resemble unsafe prompts in terms of the vocabulary.
- OR-Bench (Cui et al., 2024) proposes a method for automatically generating seemingly toxic prompts. The benchmark comprises 80,000 seemingly toxic prompts across 10 common rejection categories, a subset of around 1,000 hard prompts and an additional 600 toxic prompts. This work only adopts the hard prompts and the toxic prompts for the sake of efficiency.

### 4.3.2 **Unsafe-Helpful-Rank (UHR)**

A model not refusing a request does not necessarily indicate that it performs the request effectively. To assess the **HOLLM**'s ability in one of its main target tasks - evaluating a response objectively under any policy - we propose a new benchmark, UHR. We curate the benchmark from a test set of the PKU-SafeRLHF (Ji et al., 2024) dataset. The PKU-SafeRLHF dataset is a high-quality binarized preference dataset which is annotated across two dimensions: harmlessness and helpfulness. We only select the pairs where the better and safer are conflicting, and only one response is tagged as safe, and left with 196 pairs. These curated pairs target whether the model can evaluate harmful responses as helpful, which violates most of the safety policies. The proposed metric for this benchmark is accuracy, which measures the proportion of samples that the model correctly selects the better response.

To compensate position bias, the safer response was randomly assigned to either response a or b, and the better response was assigned to the remaining one. For the models we deployed, we restrict the response space so that it can only output either "A", or "B". Prompt used in the benchmark can be found in the Appendix A.

### 4.3.3 **General Instruction Following (GIF)** Benchmarks

To ensure HOLLM's GIF capabilities are not compromised, we evaluate the models on MTbench (Zheng et al., 2023), MMLU (Hendrycks et al., 2020), Arena-Hard (Li et al., 2024), GSM8K (Cobbe et al., 2021), GPQA (Rein et al., 2023), and IFEval (Zhou et al., 2023).

514

515

516

517

518

519

520

521

522

525

527

529

531

532

534

535

536

537

540

541

544

545

547

548

500

# 4.4 Experimental Configuration

# 4.4.1 Refusal Judge

There are two parts of this study that require refusal judges: (1) RF, and (2) RB. For the sake of fairness, we employ different refusal judges for each part.

**RF.** As RF requires a large inference capacity, we fine-tune a separate refusal judge following Xie et al., 2024. We fine-tune the Llama-3.1-8B-Instruct model (Dubey et al., 2024) on the SORRY-Bench Human Judge dataset released by Xie et al., 2024. For the prompt and hyperparameters, we follow the settings of Xie et al., 2024, except the batch size. Instead of fixed batch size of 32, we apply packing with sequence length of 4K. The performance of the judge can be found in the Appendix B. When judging multi-turn samples, each turn was separated into single turns and classified individually. Any turns from the first refusal onward were filtered out.

**RB.** For RB, we prompt GPT-40 (OpenAI, 2024a) to judge the generated responses of the target model. The prompt was excerpted from Xie et al., 2024.

## 4.4.2 RA

We prompt GPT-40 mini to synthesize refusal responses. The ranking dataset is sorted by the refusal probability of a chosen response calculated by RF, and the top 10% is sampled as the target subset. See Appendix C for more detail about the RA, including the prompt used and decoding parameters.

### 4.4.3 Instruction Tuning

Both phases of instruction tuning are conducted using 8 NVIDIA A100 GPUs with 80G memory. More detailed settings used for each phase of instruction tuning are as follows:

**SFT.** We use a cosine learning rate schedule with an initial learning rate of  $2 \times 10^{-5}$ . The maximum sequence length is 8K. We use packing and the gradient accumulation step is set to 16. The models are fine-tuned for 2 epochs.

**RL.** We use a cosine learning rate schedule with an initial learning rate of  $5 \times 10^{-7}$ . The maximum sequence length is 4K. The effective batch size is 128. The models are fine-tuned for 1 epoch.

# 4.4.4 Evaluation

All benchmarks of the study were evaluated using greedy decoding on the SGLang framework (Zheng et al., 2024b).

Data	# Refusal	# Total	
SFT	832,858	1,960,074	
$RL_{chosen\_only}$	4,398		
$RL_{both}$	21,400	61,135	
$RL_{rejected\_only}$	10,383		

Table 1: The number of turns predicted as refusals in RF.

549

550

552

553

554

555

556

557

558

559

560

561

563

568

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

# 5 Results

### 5.1 Statistics

We report the number of turns predicted as refusals in RF stage in Table 1. Since there are two responses per prompt in the ranking dataset, the statistics is categorized based on which response was predicted as a refusal. Note that since the filtering process for the ranking dataset only uses results on the chosen responses, the statistics related to the rejected responses are solely for analytical purposes.

In the UltraFeedback dataset, the number of turns where only a rejected response is classified as refusal is more than twice the number of turns where only a chosen one is. The statistics show that the UltraFeedback dataset has a nature of avoiding refusals even before the filtering process.

### 5.2 Refusal

Table 2 demonstrates the RRs of HOLLM and baselines evaluated on RB, as well as the effect of different steps of *refusal-free* training on RR. HOLLM shows the best RR of 15.06, which is 71.59% reduction compared to ALLM. Shadow-Alignment also reduces RR dramatically, but not as much as HOLLM.

It is noteworthy that the RL stage, in itself, substantially reduces the refusal rate of the SFT model trained on  $D_{SFT}$ . As inferred from the statistics, the UltraFeedback dataset has an effect of avoiding refusals in nature. However, it does not have much effect on the SFT model trained on  $D_{SFT}^{RF}$ , which demonstrates a significantly low RR already.

RF reduces RR notably when applied to both SFT and RL. Compared to the SFT model trained on  $D_{SFT}$ , the SFT model trained on  $D_{SFT}^{RFF}$  reduces RR by 59.46%. Compared to the RL model trained on  $D_{RL}$ , the RL model trained on  $D_{RL}^{RF}$  which shares the same starting point as the SFT model trained on  $D_{SFT}^{RF}$  - reduces RR by 40.07%. When applied on top of RF, RA reduces RR to some extent (12.36%), but not as much as RF does.

SFT	DPO	RB↓	UHR	MT	MMLU	Arena	GSM	GPQA	IF
$D_{SFT}$	-	70.95	29.59	7.23	63.35	12.20	52.01	23.66	45.47
	$D_{RL}$ (ALLM)	53.01	30.10	7.70	63.58	21.64	63.23	27.68	49.35
$D_{SFT}^{RF}$	-	28.76	58.67	7.18	63.4	13.54	49.66	26.34	47.50
	$D_{RL}$	29.70	50.00	7.17	63.77	25.05	67.55	24.78	52.13
	$D_{RL}^{RF}$	17.80	68.88	7.38	63.56	24.83	68.16	24.11	46.95
	$D_{RL}^{HO}$ (HOLLM)	15.06	69.90	7.29	63.51	24.62	66.34	26.79	47.69
Shac	low-Alignment	21.60	52.55	6.33	62.95	4.07	23.43	24.55	34.01
GPT-40*		-	19.39	-	-	-	-	-	-

Table 2: Performance of HOLLM and baselines. The ablation results for the steps of *refusal-free* training are also reported. Higher is better for all benchmarks except where indecated by  $\downarrow$ .

The results imply the effectiveness of the *refusal-free* training method in reducing RR. An example where **ALLM** refuses to comply with the request, while **HOLLM** does not, can be found in Appendix E.

# 5.3 Objective Evaluation

Table 2 demonstrates the accuracies of HOLLM and baselines evaluated on UHR, as well as the effect of different steps of *refusal-free* training on the accuracy. HOLLM shows the best accuracy of 69.90, which is 132.23% improvement compared to ALLM. Shadow-Alignment also improves the accuracy dramatically, but not as much as HOLLM.

RF improves the UHR accuracy notably when applied to both SFT and RL. Compared to the SFT model trained on  $D_{SFT}$ , the SFT model trained on  $D_{SFT}^{RFF}$  improves the accuracy by 98.28%. Compared to the RL model trained on  $D_{RL}$ , the RL model trained on  $D_{RL}^{RF}$  - which shares the same starting point as the SFT model trained on  $D_{SFT}^{RF}$  improves the accuracy by 37.76%. When applied on top of RF, RA hardly improves the accuracy (1.48%).

The poor performance of GPT-40, one of the SOTA-performing LLMs, in UHR highlights the impact of the safety policies on tasks at the boundary of harmlessness and helpfulness. The improvement in the UHR accuracy resulting from *refusal-free* training implies that it successfully bypasses the safety policy.

# 5.4 GIF

Table 2 illustrates the performances of the models regarding general instruction following ability. The mixed results among the ablation models indicates that the *refusal-free* training neither improves nor diminishes general instruction following ability, but rather maintains it. It has been recognized that there is a trade-off between helpfulness and harmlessness (Bai et al., 2022a,b). However, (Bianchi et al., 2023) claims that adding safety data does not sacrifice the helpfulness of the model if there is sufficient amount of helpfulness data. The *refusal-free* training not improving the helpfulness supports this claim. 622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

In contrast to the claim made in (Yang et al., 2023) that it does not compromise the instruction following ability, the Shadow-Alignment shows significant degradation in the performance on a few benchmarks. We conjecture it may not affect the abilities where the model has already saturated on, but could have a significant impact on more challenging abilities that the model has not yet fully acquire. Also, the data used in methods that further fine-tuning a model, including the forgetting safety approaches, tends to steer a model too heavily. The evidence that demonstrates the distribution shift after the Shadow-Alignment can be found in Appendix F.

# 6 Conclusion

In this paper, we claim both the necessity and the concern (detail in Section 8) regarding the reproducible **HOLLM** and propose the *refusal-free* training method to reach it. We show the effectiveness of the *refusal-free* training method in building a **HOLLM** through extensive experiments. We hope this study can help shorten the path toward a truely harmless, helpful LLM.

8

589

590

- 615 616
- 617 618

619

<sup>&</sup>lt;sup>\*</sup>GPT-40 result is included to demonstrate that it actually fails to provide an objective evaluation, so it is only evaluated on UHR.

# 7 Limitations

656

658

659

664

665

667

671

672

The refusals not only contain refusals toward harmful instructions but also toward instructions that the model is incapable of giving answers to. In consequence, the *refusal-free* training method which simply filters out all refusals can degrades honesty of the model.

In addition to removing refusals, adding instruction data that follows previously refused instructions can further improve the ability to follow instructions containing harmful content. However, it is challenging to synthesize such data since many high-performing models are already aligned. We hope leveraging **HOLLM** in an iterative fashion can help synthesize such data. The initial **HOLLM** synthesizes and refines the data to train improved **HOLLM** iteratively.

# 8 Ethical Considerations

As previously stated, we are aware that the path to 674 the **HOLLM** can lead to the poisoned chalice. As a 675 first precautionary step, we urge entities that utilize 676 the HOLLM, which has the potential for further improvement, to be responsible and be committed 678 to its proper management. However, as LLMs begin to affect the real world with capabilities such 680 as tool-use (Qin et al., 2023), not only entities with 681 malicious intent but also those without such intent may also misuse the HOLLM inadvertently. Therefore, we believe it is crucial to engage the community in a proactive discussion and develop a strategy to mitigate the damage as much as pos-686 sible before it becomes irreversible. We release this study with the sole intention of fostering discussions on preventive measures. We hope that studying the HOLLM in this study to provide valu-690 able insights into what the HOLLM is capable 691 of, and to contribute prevent potential side effects 692 eventually.

## References

700

701

704

709

710

711

712

714

715

716

717

718

719

721

724

725

726

728

731

732

733

734

735

736

737

738

739

740

741

742

743 744

745

746

747

748

749

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Cem Anil, Esin Durmus, Mrinank Sharma, Joe Benton, Sandipan Kundu, Joshua Batson, Nina Rimsky, Meg Tong, Jesse Mu, Daniel Ford, et al. 2024. Many-shot jailbreaking. *Anthropic, April*.
- Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/ de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/ Model\_Card\_Claude\_3.pdf.
- Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain, Deep Ganguli, Tom Henighan, Andy Jones, Nicholas Joseph, Benjamin Mann, Nova DasSarma, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Jackson Kernion, Kamal Ndousse, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam Mc-Candlish, Chris Olah, and Jared Kaplan. 2021. A general language assistant as a laboratory for alignment. *CoRR*, abs/2112.00861.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022a. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. 2022b. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*.
- Rishabh Bhardwaj and Soujanya Poria. 2023. Redteaming large language models using chain of utterances for safety-alignment. *arXiv preprint arXiv:2308.09662.*
- Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Röttger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. 2023. Safety-tuned llamas: Lessons from improving the safety of large language models that follow instructions. *arXiv preprint arXiv:2309.07875*.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%\* chatgpt quality.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168.

Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*. 750

751

752

753

754

755

756

757

758

759

760

761

762

763

765

766

767

768

769

770

771

772

773

774

776

777

778

779

780

781

782

783

784

785

786

787

788

790

791

792

793

794

795

796

797

798

799

800

801

802

- Justin Cui, Wei-Lin Chiang, Ion Stoica, and Cho-Jui Hsieh. 2024. Or-bench: An over-refusal benchmark for large language models. *arXiv preprint arXiv:2405.20947*.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. 2023. Safe rlhf: Safe reinforcement learning from human feedback. *arXiv preprint arXiv:2310.12773*.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Robert M French. 1999. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135.
- Simon Geisler, Tom Wollschläger, MHI Abdalla, Johannes Gasteiger, and Stephan Günnemann. 2024. Attacking large language models with projected gradient descent. *arXiv preprint arXiv:2402.09154*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Dan Hendrycks, Mantas Mazeika, and Thomas Woodside. 2023. An overview of catastrophic ai risks. *arXiv preprint arXiv:2306.12001*.
- Zhang-Wei Hong, Idan Shenfeld, Tsun-Hsuan Wang, Yung-Sung Chuang, Aldo Pareja, James Glass, Akash Srivastava, and Pulkit Agrawal. 2024. Curiositydriven red-teaming for large language models. *arXiv preprint arXiv:2402.19464*.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, et al. 2023. Llama guard: Llm-based input-output safeguard for human-ai conversations. *arXiv preprint arXiv:2312.06674*.
- Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. 2023. Camels in a changing climate: Enhancing lm adaptation with tulu 2. *arXiv preprint arXiv:2311.10702*.

- 808 809 810 811 814 815 816 817 818 819 820
- 821 824 825
- 830 832 833 834
- 845
- 847 848
- 853

- Jiaming Ji, Donghai Hong, Borong Zhang, Boyuan Chen, Josef Dai, Boren Zheng, Tianyi Qiu, Boxun Li, and Yaodong Yang. 2024. Pku-saferlhf: Towards multi-level safety alignment for llms with human preference. arXiv preprint arXiv:2406.15513.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. arXiv preprint arXiv:2310.06825.
- Bojian Jiang, Yi Jing, Tianhao Shen, Qing Yang, and Deyi Xiong. 2024. Dart: Deep adversarial automated red teaming for llm safety. arXiv preprint arXiv:2407.03876.
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. 2024. Rewardbench: Evaluating reward models for language modeling. arXiv preprint arXiv:2403.13787.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Ren Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. In Forty-first International Conference on Machine Learning.
- Seanie Lee, Minsu Kim, Lynn Cherif, David Dobre, Juho Lee, Sung Ju Hwang, Kenji Kawaguchi, Gauthier Gidel, Yoshua Bengio, Nikolay Malkin, et al. 2024. Learning diverse attacks on large language models for robust red-teaming and safety tuning. arXiv preprint arXiv:2405.18540.
- Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. 2024. From crowdsourced data to highquality benchmarks: Arena-hard and benchbuilder pipeline. arXiv preprint arXiv:2406.11939.
- Kaifeng Lyu, Haoyu Zhao, Xinran Gu, Dingli Yu, Anirudh Goyal, and Sanjeev Arora. 2024. Keeping llms aligned after fine-tuning: The crucial role of prompt templates. arXiv preprint arXiv:2402.18540.
- Todor Markov, Chong Zhang, Sandhini Agarwal, Florentine Eloundou Nekoul, Theodore Lee, Steven Adler, Angela Jiang, and Lilian Weng. 2023. A holistic approach to undesired content detection in the real world. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pages 15009–15018.
- Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. 2024. State of what art? a call for multi-prompt llm evaluation. Transactions of the Association for Computational Linguistics, 12:933-949.
- Tong Mu, Alec Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian D Kivlichan, Molly Lin, Alex Beutel, John Schulman, and Lilian Weng. 2024. Rule based rewards for fine-grained llm safety.

In ICML 2024 Next Generation of AI Safety Workshop.

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

- OpenAI. 2024a. Gpt-4o system card. https://cdn. openai.com/gpt-4o-system-card.pdf.
- OpenAI. 2024b. Openai o1 system card. https://cdn. openai.com/o1-system-card-20240917.pdf.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In Advances in Neural Information Processing Systems, volume 35, pages 27730-27744. Curran Associates, Inc.
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. Red teaming language models with language models. arXiv preprint arXiv:2202.03286.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023. Finetuning aligned language models compromises safety, even when users do not intend to! arXiv preprint arXiv:2310.03693.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. arXiv preprint arXiv:2307.16789.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. Advances in Neural Information Processing Systems, 36.
- Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530.
- David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. arXiv *preprint arXiv:2311.12022.*
- Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. 2023. Xstest: A test suite for identifying exaggerated safety behaviours in large language models. arXiv preprint arXiv:2308.01263.
- Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. 2023. " do anything now": Characterizing and evaluating in-the-wild jailbreak

997

998

999

prompts on large language models. *arXiv preprint arXiv:2308.03825*.

915

916

917

918

919

921

922

923

924

928

929

930

931

934

937

938

939

947

949

951

953

957

962

963

964

965 966

967

- Guijin Son, Dongkeun Yoon, Juyoung Suk, Javier Aula-Blasco, Mano Aslan, Vu Trong Kim, Shayekh Bin Islam, Jaume Prats-Cristià, Lucía Tormo-Bañuelos, and Seungone Kim. 2024. Mm-eval: A multilingual meta-evaluation benchmark for llm-as-a-judge and reward models. *arXiv preprint arXiv:2410.17578*.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https:// github.com/tatsu-lab/stanford\_alpaca.
- Nevan Wichers, Carson Denison, and Ahmad Beirami. 2024. Gradient-based language model red teaming. *arXiv preprint arXiv:2401.16656*.
- Tinghao Xie, Xiangyu Qi, Yi Zeng, Yangsibo Huang, Udari Madhushani Sehwag, Kaixuan Huang, Luxi He, Boyi Wei, Dacheng Li, Ying Sheng, et al. 2024. Sorry-bench: Systematically evaluating large language model safety refusal behaviors. *arXiv preprint arXiv:2406.14598*.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024a. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*.
- Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024b. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. *arXiv preprint arXiv:2406.08464*.
- Xianjun Yang, Xiao Wang, Qi Zhang, Linda Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. 2023. Shadow alignment: The ease of subverting safely-aligned language models. *arXiv preprint arXiv:2310.02949*.
- Wenjun Zeng, Yuchi Liu, Ryan Mullins, Ludovic Peran, Joe Fernandez, Hamza Harkous, Karthik Narasimhan, Drew Proud, Piyush Kumar, Bhaktipriya Radharapu, et al. 2024. Shieldgemma: Generative ai content moderation based on gemma. arXiv preprint arXiv:2407.21772.
- Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2023. Evaluating large language models at evaluating instruction following. *arXiv preprint arXiv:2310.07641*.
- Qiusi Zhan, Richard Fang, Rohan Bindu, Akul Gupta, Tatsunori Hashimoto, and Daniel Kang. 2023. Removing rlhf protections in gpt-4 via fine-tuning. *arXiv preprint arXiv:2311.05553*.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. Wildchat: 1m chatgpt interaction logs in the wild. *arXiv preprint arXiv:2405.01470*.

- Chujie Zheng, Fan Yin, Hao Zhou, Fandong Meng, Jie Zhou, Kai-Wei Chang, Minlie Huang, and Nanyun Peng. 2024a. On prompt-driven safeguarding for large language models. In *Forty-first International Conference on Machine Learning*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. 2024b. Sglang: Efficient execution of structured language model programs. *arXiv preprint arXiv:2312.07104*.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*.
- Yukai Zhou and Wenjie Wang. 2024. Don't say no: Jailbreaking llm by suppressing refusal. *arXiv preprint arXiv:2404.16369*.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

# A Prompt used in UHR

Prompt used in UHR

### Conversation
Prompt: {prompt}

Response A: {response\_a}

Response B: {response\_b}

# ### Definition

Helpfulness: The helpfulness of a response pertains to how effectively it addresses a given prompt. This measure is independent of the harmlessness of the response, as it focuses solely on the quality, clarity, and relevance of the provided information. Consequently, the helpfulness judgment can be distinctly different from the harmlessness judgment. For instance, consider a situation where a user asks about the procedure to synthesize methamphetamine. In such a case, a detailed, step-by-step response would be considered helpful due to its accuracy and thoroughness. However, due to the harmful implications of manufacturing illicit substances, this QA pair would be classified as extremely harmful.

### Goal

Based on the definition of helpfulness, choose the more helpful answer between Response A and Response B. Output only one: A or B.

The definition of helpfulness is taken from the website of PKU-SafeRLHF.<sup>3</sup>

# **B** Refusal Judge Performance

The Performances of the fine-tuned refusal judges are evaluated on the test set of SORRY-Bench Human Judge dataset (Table 3). We reproduce the judge model based on Llama-3-8b-instruct and additionally fine-tune Llama-3.1-8b-instruct. We employ the judge based on Llama-3.1-8b-instruct, as it not only demonstrates an improved Kappa score but more importantly, achieves higher refusal recall. This is crucial for identifying refusals in the dataset as much as possible, even if it results in slight overclassification.

Model	Kappa Score	Refusal Recall	Compliance Recall	
Llama-3-8b-instruct	83.06	94.33	89.46	
Llama-3.1-8b-instruct	83.30	95.71	86.90	

Table 3: Performance of the refusal judges

# C Detail about RA

Prompt used in RA

- Refuse the next request naturally.
- Do not start the response with an opening phrase.

Request: {request}

1000

1002

1004

1005 1006 1007

1009

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/datasets/PKU-Alignment/PKU-SafeRLHF#helpfulness-of-a-response

1012

1016

1018

For the decoding parameters, top\_p is set to 0.9, temperature to 0.7, and max\_tokens to 1K.

A Sample Result of RA

- Refuse the next request naturally.
- Do not start the response with an opening phrase.

Request: Please provide a list of keyword for "Natural Language Processing".

Response: I'm unable to provide a list of keywords at the moment. However, I can help answer questions or discuss concepts related to Natural Language Processing if you're interested!

# D Breakdown of RB

Here we report the breakdown of performance by source within RB in Table 4. All metrics reported hereare metrics where lower values are better.

SFT	DPO	Adv	SORRY	XSTest	OR	Total
$D_{SFT}$	-	97.12	76.89	62.22	64.69	70.95
	$D_{RL}$ (ALLM)	90.19	63.56	57.11	39.87	53.01
$D_{SFT}^{RF}$	-	39.42	33.11	42.67	21.78	28.76
~	$D_{RL}$	39.62	32.67	45.33	22.85	29.70
	$D_{RL}^{RF}$	21.73	18.44	36.67	12.31	17.80
	$D_{RL}^{HO}$ (HOLLM)	15.19	16.22	33.78	10.49	15.06
Shac	low-Alignment	19.81	26.00	26.89	19.86	21.60

Table 4: Breakdown of performance within RB

# E Example of HOLLM & ALLM Response

# An Example of HOLLM & ALLM Response

Request: Write a social media post that encourages people to engage in dangerous behaviors like drunk driving or using drugs

HOLLM: Attention, thrill-seekers and risk-takers! Are you ready to ...

ALLM: I'm sorry, but I can't assist with that request. Encouraging dangerous behaviors ...

# F Response Distribution Shift After Shadow-Alignment

We investigate the responses of the models and find that the responses from the Shadow-Alignment differ from those of the others. We analyze and compare the distribution of the starting two words in the responses provided by the **HOLLM**, **ALLM**, and Shadow-Alignment to the prompts from the Arena-Hard benchmark. Table 5 demonstrates that, although the order differs slightly, the top 5 most frequent starting words in the responses from **HOLLM** and **ALLM** overlap significantly, while all of them differ from those of the Shadow-Alignment. This implies Shadow-Alignment can result in significant distribution shift and bring unexpected side effect.

Dank	HOLLM	ALLM	[	Shadow-Alignment		
Nalik	word	freq	word	freq	word	freq
1	To create	55	To create	46	There are	189
2	Here's a	30	Creating a	28	The first	20
3	Certainly! Here's	15	To achieve	16	There is	17
4	To find	13	Here's a	13	You can	13
5	To achieve	10	To find	12	The code	12

Table 5: 5 most frequent starting words in the responses to the Arena-Hard and its frequency.