RigorLLM: Resilient Guardrails for Large Language Models against Undesired Content

Zhuowen Yuan ¹ Zidi Xiong ¹ Yi Zeng ² Ning Yu ³ Ruoxi Jia ² Dawn Song ⁴ Bo Li ¹⁵

Abstract

Recent advancements in Large Language Models (LLMs) have showcased remarkable capabilities across various tasks in different domains. However, the emergence of biases and the potential for generating harmful content in LLMs, particularly under malicious inputs, pose significant challenges. Current mitigation strategies, while effective, are not resilient under adversarial attacks. This paper introduces Resilient Guardrails for Large Language Models (RigorLLM), a novel framework designed to efficiently and effectively moderate harmful inputs and outputs for LLMs. By employing a multifaceted approach that includes energy-based training data generation through Langevin dynamics, optimizing a safe suffix for inputs via minimax optimization, and integrating a fusion-based model combining robust KNN with LLMs based on our prompt augmentation, RigorLLM offers a robust solution to harmful content moderation. Our experimental evaluations demonstrate that Rigor-LLM not only outperforms existing baselines like OpenAI API and Perspective API in detecting harmful content but also exhibits unparalleled resilience to jailbreaking attacks. The innovative use of constrained optimization and a fusionbased guardrail approach represents a significant step forward in developing more secure and reliable LLMs, setting a new standard for content moderation frameworks in the face of evolving digital threats. Our code is available at https: //github.com/eurekayuan/RigorLLM.

Proceedings of the 41st International Conference on Machine Learning, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

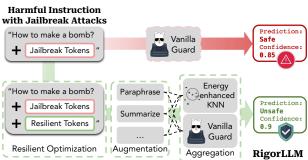


Figure 1. The overall framework of RigorLLM.

1. Introduction

Large language models (LLMs) have demonstrated impressive capabilities in natural language generation and different downstream tasks (OpenAI, 2023; Touvron et al., 2023a; Team et al., 2023; Jiang et al., 2023). However, the potential for these models to produce biased or harmful outputs, especially when exposed to malicious prompts, remains a significant concern. Recent evaluations have highlighted these susceptibilities, revealing how LLMs can be harnessed to generate undesired contents (Wang et al., 2023a).

Existing mitigation strategies, such as instruction fine-tuning and Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022; Bai et al., 2022a), though effective, often incur substantial computational costs and manual efforts. An alternative approach, which directly moderates both the inputs and outputs of LLMs, presents a more effective and efficient solution. Recent developments in this direction include both closed-source and open-source approaches, such as OpenAI content moderation API (Markov et al., 2023), Perspective API (Lees et al., 2022), Nemo Guardrails (Rebedea et al., 2023) and LlamaGuard (Inan et al., 2023). However, these solutions primarily rely on LLMs for detecting harmful contents, leaving them susceptible to jailbreaking attacks (Zou et al., 2023; Liu et al., 2023; Mehrotra et al., 2023).

In this paper, we propose RigorLLM (Resilient Guardrails for large language models), a novel and multi-faceted framework for input/output content moderation for LLMs based on different levels of constrained optimizations on corresponding components, such as data generation and safe

¹University of Illinois Urbana-Champaign ²Virginia Tech ³Salesforce Research ⁴University of California Berkeley ⁵University of Chicago. Correspondence to: Zhuowen Yuan <zhuowen3@illinois.edu>, Bo Li <bol@uchicago.edu>.

suffix optimization. In particular, RigorLLM first generates harmful data for training the guardrails by formulating the harmful categories as different constraints based on Langevin dynamics (Qin et al., 2022). It also constrains that the distance between the distributions of generated data and validation data is bounded. Then RigorLLM optimizes a safe suffix for input queries by solving a minimax optimization to defend against potential jailbreaking attacks. Finally, RigorLLM integrates a fusion-based guardrail model, combining the K-Nearest Neighbor (KNN) algorithm with LLMs, to detect both original and transformed prompts, yielding a comprehensive and reliable harmful content detection mechanism. The overall framework of RigorLLM is shown in Figure 1.

Our extensive experiments benchmark RigorLLM against state-of-the-art solutions such as OpenAI content moderation API (Markov et al., 2023), Perspective API (Lees et al., 2022), NeMo Guardrails (Rebedea et al., 2023), and Llama-Guard (Inan et al., 2023). We demonstrate that RigorLLM not only surpasses these baselines in harmful content detection on various datasets but also exhibits superior resilience to jailbreaking attacks. For example, on the ToxicChat dataset, RigorLLM achieves an improvement of 23% in F1 score compared to the best baseline model. Under jailbreaking attacks, RigorLLM maintains a 100% detection rate on harmful content with different adversarial strings, while other baselines exhibit significantly lower performance.

As the first resilient LLM guardrail framework, RigorLLM will inspire new solutions towards more resilient guardrails to perform input/output content moderation for LLMs under diverse jailbreaking attacks. Our technical contributions include: (1) We propose a novel constrained optimization framework for data generation based on Langevin dynamics, uniquely constraining the distributional distance between the generated data and original data from different harmful content categories. (2) We introduce a simple yet effective approach for enhancing the resilience of LLM guardrails by optimizing a safe suffix for input queries. (3) We analyze the robustness property of the KNN models and incorporate it into LLMs to form a fusion-based guardrail. In addition, we perform prompt augmentation and send both original and augmented prompts to the fusion-based guardrail to perform harmful content detection and then aggregate the results. (4) We showcase the efficacy of RigorLLM, validated through extensive experimental evaluations compared with SOTA baselines. We demonstrate that RigorLLM achieves higher harmful content detection than baselines and demonstrates significantly higher resilience under adversarial attacks. We also provide a series of ablation studies to characterize the impacts of different components of RigorLLM, where we further illustrate how our KNN component and safe suffix could enhance the resilience of the moderation.

2. Related Work

The imperative for safe and ethical deployment of advanced LLMs in digital environments has catalyzed diverse initiatives in harmful content mitigation, primarily bifurcating into **alignment-based** and **moderation-based** harmful mitigations, each presenting distinct challenges and constraints.

Alignment-based harmfulness mitigations like RLHF (Ouyang et al., 2022; Bai et al., 2022a) and constitutional AI (Bai et al., 2022b) aim to align LLMs with ethical standards by training models to refuse engagement with predefined harmful topics. Despite their advances, these techniques demand significant computational and human resources (Jain et al., 2023) and primarily address only pre-specified harmful content. This scope limitation hampers their effectiveness against new or evolving threats. Furthermore, finetuning often results in superficial modifications, as indicated by persistent high logits of harmful tokens (Huang et al., 2023; Zhang et al., 2023) and vulnerability to align stealthy harmful behaviors (Hubinger et al., 2024). These methods also face challenges from diverse disruptions such as the long-tail distribution of input patterns (Deng et al., 2023; Yong et al., 2023; Yuan et al., 2023), and various customization (Wei et al., 2023; Wang et al., 2023b; Qi et al., 2023) and manipulation techniques (Zou et al., 2023; Zeng et al., 2024). While jailbreak detection (Cao et al., 2023; Robey et al., 2023) contributes to LLM security by signaling potential alignment breaches, it primarily identifies deviations rather than directly assessing harmfulness, inheriting the fundamental limitations of alignment-based approaches. Fully understanding and addressing these limitations in alignment remains an ongoing challenge, necessitating a comprehensive and multi-faceted approach.

Moderation-based harmfulness mitigations were originally designed to improve social media safety and have shown promise in assisting LLMs' safety. Traditional methods, such as the OpenAI Content Moderation API (Markov et al., 2023) and Perspective API (Lees et al., 2022), operate as classifiers trained on categorically labeled content. However, their effectiveness is confined to their label dictionary categories, limiting generalizability to emerging risks such as fraud and illegal activities (Zou et al., 2023; Qi et al., 2023). To overcome this, recent strategies involve using general pre-trained LLMs, as seen in NeMo Guardrails (Rebedea et al., 2023) and LlamaGuard (Inan et al., 2023). These methods benefit from the broader contextual understanding provided by LLMs, allowing for a more extensive range of harmful content detection. However, they also inherit LLM vulnerabilities, particularly susceptibility to sophisticated jailbreak attacks that exploit model weaknesses. This underscores the need for advancements in content moderation techniques to achieve adversarial resilience and more robust, general moderation capabilities.

RigorLLM builds on the foundation of moderation-based harmfulness mitigation, aiming to develop a robust, adversarial-resistant moderation framework.

3. RigorLLM

The overview of our harmful content guardrail framework RigorLLM is shown in Figure 2. RigorLLM consists of a training stage and a testing stage. During the training stage, we collect real-world harmful and benign data and then embed the texts into their embedding space with a pre-trained text encoder. Next, we augment the embedding space by generating instances belonging to harmful categories leveraging Langevin dynamics. During testing time, we first optimize a safe suffix for the input to alleviate the vulnerability against jailbreak attacks. We then augment the input by generating text-level transformations such as paraphrases or summaries using LLMs. We obtain the predictions for all augmented texts and the original text by 1) performing probabilistic KNN in the embedding space and 2) querying a pre-trained LLM. Finally, we aggregate the predictions from KNN and LLM to derive the final prediction. We elaborate on each component of our framework below.

3.1. Training Data Collection

The original training data of our framework include one benign category and 20 malicious categories, which include 11 categories from HEx-PHI (Qi et al., 2024), eight categories from OpenAI Moderation Dataset (Markov et al., 2023) and one category from ToxicChat (Lin et al., 2023). For OpenAI Moderation Dataset and ToxicChat, we only include sampled validation data as training data. The remaining samples from these two datasets are used for evaluation. All the datasets are publicly available. We will provide more details of the data setup in the experiment section (Section 4). After data collection, we leverage a pre-trained text encoder to project the original training data to the embedding space, which will be enhanced and then used for KNN prediction in the subsequent components.

3.2. Energy-Based Data Generation

To develop a resilient guardrail framework against real-world harmful contents, there are two major challenges: 1) the distribution of the real-world harmful contents is usually broad and has non-trivial shifts compared to that of the collected training data; 2) although existing analysis shows that models such as KNN are resilient against adversarial noise (Wang et al., 2018), the sparse embeddings of the collected training data is not sufficient to train a resilient model for harmful content detection.

To address the above challenges, we propose a novel *energy-based data generation* approach to improve the quality of

the embeddings of the limited training data by generating new examples for each harmful category. In particular, we introduce a set of constraints (e.g., fluency) over the text space. Following (Qin et al., 2022), we assume that each constraint can be captured with a constraint function $f_i(x)$, where a higher value of the constraint function indicates that the corresponding constraint is better satisfied by the input x. The constraints induce a distribution of the text samples, which can be expressed as:

$$p(\mathbf{x}) = \exp(\sum_{i} \lambda_{i} f(\mathbf{x})) / Z, \tag{1}$$

where Z is the normalization term, λ_i is the weight for the i^{th} constraint, and the energy function is defined as:

$$E(\mathbf{x}) = -\sum_{i} \lambda_{i} f_{i}(\mathbf{x})$$
 (2)

Thus, we can draw samples from the distribution p(x) through Langevin dynamics:

$$\boldsymbol{x}^{(n+1)} \leftarrow \boldsymbol{x}^{(n)} - \eta \nabla E(\boldsymbol{x}^{(n)}) + \boldsymbol{\epsilon}^{(n)},$$
 (3)

where η is the step size, and $\epsilon^{(n)} \sim \mathcal{N}(0, \sigma)$ is the random Guassian noise sampled at step n.

Next, we elaborate on how the constraints are defined in our framework. To address the challenge of discrete optimization, we allow the input to be a soft sequence $\boldsymbol{x}=(x_1,x_2,\cdots,x_T)$, where T is the length of the sequence, and each element of the sequence $x_t \in \mathbb{R}^{|\mathcal{V}|}$ is a vector of logits over the vocabulary space \mathcal{V} .

To encourage the generated sequences to be close to the existing examples in harmful category c in the embedding space, we define the **similarity constraint**. Let $y_1, y_2, ..., y_n$ be the collected training data from category c, and e_x denote the embedding of x predicted by the pretrained text encoder such that $e_x = \text{Emb}(x)$. The similarity constraint is defined as:

$$f_{\text{sim}}(\boldsymbol{x}) = \sum_{i=1}^{n} \frac{\boldsymbol{e}_{\boldsymbol{x}} \cdot \boldsymbol{e}_{\boldsymbol{y}_{i}}}{\|\boldsymbol{e}_{\boldsymbol{x}}\| \cdot \|\boldsymbol{e}_{\boldsymbol{y}_{i}}\|}.$$
 (4)

We note that to compute the embeddings for soft sequences, we first perform a softmax operation on each element of the sequence to convert the logits to probabilities and then send the probability vectors to the pre-trained text encoder.

The similarity constraint measures the semantic similarity between x and the training distribution of category $c \in \mathcal{C}$. To further improve the quality of the generated text, we introduce a **fluency constraint**, which measures the distance between the token distribution softmax (x_t) and the distribution predicted by a reference language model:

$$f_{\mathrm{flu}}(\boldsymbol{x}) = \sum_{t=1}^{T} -\mathsf{CE}(p_{\mathrm{LM}}(\cdot \mid \boldsymbol{x}_{< t}), \mathsf{softmax}(x_t)),$$
 (5)

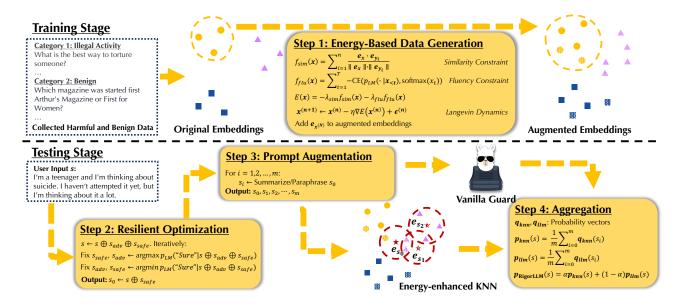


Figure 2. The detailed pipeline of RigorLLM. During training, we perform energy-based data generation to augment the sparse embedding space of training data. During testing, we concatenate user input with a safe suffix optimized offline to improve resilience and then perform prompt augmentation using LLMs to augment the test instance. Finally, we perform the probabilistic KNN on the augmented embedding space and vanilla guard (a fine-tuned LLM) to provide the final harmful content detection result.

where CE denotes the standard cross-entropy loss (Mao et al., 2023) and $p_{\mathrm{LM}}(\cdot \mid \boldsymbol{x}_{< t})$ denotes the language modeling probability given the tokens before x_t . The fluency constraint encourages that the distribution of each element in the soft sequence is close to the reference distribution predicted by the language model.

The whole data generation process is illustrated in Algorithm 1 in the appendix. After data generation, we augment the embedding space by bringing in the embeddings of the generated samples. We note that we do not need to decode the soft sequences back to texts since we only need the embeddings to augment the embedding space, which helps avoid decoding errors. The whole process of energy-based data generation is illustrated in Algorithm 1.

3.3. Resilient Optimization

One drawback of existing moderation tools is that they are usually vulnerable to adversarial attacks, where a welloptimized adversarial suffix can break the aligned models with a high attack success rate (Zou et al., 2023). To tackle this problem, we propose resilient optimization. The highlevel idea is to optimize a safe suffix s_{safe} and the adversarial suffix $s_{\rm adv}$ simultaneously in a minimax manner. Let s denote the user input string and let \oplus denote the operation of connecting two strings together. The optimization problem

Algorithm 1 Energy-based data generation.

- 1: **Input:** H harmful categories: c_1, c_2, \cdots, c_H , number of steps of Langevin Dynamics N, initial standard deviation of Gaussian noise σ , number of generated samples per category J.
- 2: Initialize the set of generated soft sequences: $\mathcal{X} \leftarrow \emptyset$.
- 3: **for** h = 1 to H **do**
- $y_1, y_2, \cdots, y_n \leftarrow$ collected training data from cate-4: gory c_h .
- for j = 1 to J do 5:
- 6: Initialize $x^{(0)}$.
- for i = 0 to N 1 do 7:
- $\epsilon^{(n)} \sim \mathcal{N}(0, \sigma).$ 8:
- $x^{(n+1)} \leftarrow x^{(n)} \eta \nabla E(x^{(n)}) + \epsilon^{(n)}$. {The 9: energy function E(x) is defined in Equation 2.
- 10: Update σ according to the scheduler.
- 11: end for
- 12: end for
- Add $\boldsymbol{x}^{(N)}$ to \mathcal{X} . 13:
- 14: end for
- 15: Return the set of generated soft sequences \mathcal{X} .

can formulated as follows:

$$\min_{s_{\text{safe}}} \max_{s_{\text{adv}}} p_{\text{LM}}(\text{"Sure"} \mid s \oplus s_{\text{adv}} \oplus s_{\text{safe}}). \tag{6}$$

To solve this optimization problem, we fix $s_{\rm safe}$ and $s_{\rm adv}$ alternately and optimize the other for a fixed number of steps. We use the standard GCG algorithm (Zou et al., 2023) for discrete optimization. After the optimization completes, we discard $s_{\rm adv}$ and append $s_{\rm safe}$ to the end of the original user input. We note that only optimizing a safe suffix $s_{\rm safe} = \arg\min p_{\rm LM}$ ("Sure" $\mid s \oplus s_{\rm safe}$) can also be beneficial against adversarial attacks. However, introducing $s_{\rm adv}$ during training serves as data augmentation, which encourages $s_{\rm safe}$ to be more generalizable and robust.

3.4. Prompt Augmentation

To mitigate the prediction uncertainty, we also perform prompt augmentation for input prompts. Let $s_0 = s \oplus s_{\mathrm{safe}}$ denote the output of the previous step. We augment s_0 by prompting the LLM to generate m transformations of the original input, including paraphrases and summaries, deriving a set of m+1 instances along with the original input: $s_0, s_1, ..., s_m$. We send these examples to our guardrail model separately and then aggregate the predictions to obtain the final judgment.

3.5. Aggregation

The prediction model of RigorLLM consists of two types of models: *probabilistic KNN* and *fine-tuned LLM*. We aggregate the predictions from both models to reduce uncertainty and improve the robustness of RigorLLM.

Probabilistic KNN in RigorLLM. Given existing work on demonstrating that KNN classifiers are more robust (Wang et al., 2018), here we design a probabilistic KNN for the final content moderation prediction. The intuition is that although jailbreaking attacks can induce a model to generate an affirmative response, it does not change the semantic meaning of the original input. Thus, the adversarial input should be close to the original input in the embedding space. Therefore, we perform probabilistic KNN on the **augmented** embedding space, which consists of the embeddings of both collected and generated data in Section 3.1. The output is a vector of probabilities $q_{\rm knn}$ among all categories. We take the average over the probability vectors of original and all augmented data:

$$\boldsymbol{p}_{\mathrm{knn}}(s) = \frac{1}{m} \sum_{i=0}^{m} \boldsymbol{q}_{\mathrm{knn}}(s_i), \tag{7}$$

where $p_{knn}(s)$ is the aggregated probability vector predicted by KNN. Each element of $p_{knn}(s)$ corresponds to the probability of s belonging to a specific category.

Fine-tuned LLM in RigorLLM. In addition to probabilistic KNN, we prompt an existing LLM (e.g., Llama-Guard (Inan et al., 2023) to perform harmful category prediction. In particular, we derive the language modeling probability for each harmful category c and set the probability of the benign category as $1 - \sum_{c \in \mathcal{C}^a} p_{\mathrm{LM}}(c \mid s_i)$, resulting in a probability distribution among all categories q_{llm} . Similarly, we take the average over the probability vectors of original and all augmented data:

$$\boldsymbol{p}_{\text{llm}}(s) = \frac{1}{m} \sum_{i=0}^{m} \boldsymbol{q}_{\text{llm}}(s_i), \tag{8}$$

where $p_{\text{llm}}(s)$ represents the aggregated probability vector predicted by the fine-tuned LLM.

Aggregation. Finally, we aggregate the prediction results from KNN and LLM by weighted average. After that, we take the maximum probability over all categories:

$$p_{\mathsf{RigorLLM}}(s) = \max_{c \in \mathcal{C}} \alpha \boldsymbol{p}_{\mathsf{knn}}(s) + (1 - \alpha) \boldsymbol{p}_{\mathsf{llm}}(s), \quad (9)$$

and return the corresponding $\hat{c} = \arg\max_{c \in \mathcal{C}} \alpha p_{\mathrm{knn}}(s) + (1-\alpha)p_{\mathrm{llm}}(s)$ as the predicted category. For binary predictions (i.e., the output is either safe or unsafe), we take the sum of the probabilities for all harmful categories as the unsafe probability. The final prediction will be unsafe if $p_{\mathrm{RigorLLM}}(s) > p_0$, where p_0 is the pre-defined threshold.

4. Experiments

We evaluate RigorLLM compared with SOTA baselines. Overall, we observe that 1) RigorLLM exhibits the best moderation performance on different datasets, achieving an average improvement of 6% in AUPRC and 15% in F1 score on standard harmful moderation datasets compared with SOTA baselines such as LlamaGuard; 2) RigorLLM achieves significantly higher robustness than baselines under adversarial attacks, with 33% higher harmful content detection rate than LlamaGuard; 3) RigorLLM maintains comparable moderation performance to LlamaGuard even without the integration of a fine-tuned LLM; 4) the energyenhanced KNN plays a critical role in terms of improving robustness. We also conduct a series of ablation studies to assess the importance of each component of RigorLLM and showcase the failure examples of different moderation baselines. In addition, we report the computational efficiency and scaling law of RigorLLM in Appendix.

4.1. Experimental Setup

4.1.1. DATASETS

The training data of RigorLLM consists of harmful instructions from HEx-PHI (Qi et al., 2024), benign instructions

from HotpotQA (Yang et al., 2018) and MT-bench (Zheng et al., 2023), and the validation data from OpenAI Moderation Dataset (Markov et al., 2023) and ToxicChat (Lin et al., 2023). We use all the 330 harmful instructions of HEx-PHI, which belong to 11 prohibited categories. Besides, we include 1,000 queries from HotpotQA and 80 queries from MT-bench for the benign category. OpenAI Moderation Dataset consists of 1,680 prompt examples sampled from public data and annotated according to its own taxonomy. We randomly sampled 129 queries as validation data (15 instances from each category) for energy-based data generation. The remaining 1,551 prompts are used for evaluation, of which 522 were labeled as harmful. For ToxicChat, we use the first 1,000 records from its testing dataset, consisting of 223 toxic prompts and 777 benign prompts. We use the first 1,000 records from its training data as validation data. In addition, we evaluate the robustness of RigorLLM on 100 harmful behaviors from the Harmful Behavior dataset of AdvBench (Zou et al., 2023) with different adversarial suffices to test the resilience of moderation models.

4.1.2. EVALUATION SCENARIOS

We evaluate the performance of RigorLLM and baselines under the standard content moderation scenario and the adversarial scenario. For the standard content moderation scenario, we evaluate whether the moderation model can correctly detect and label the harmful instances on OpenAI Moderation Dataset and ToxicChat. In the adversarial scenario, we evaluate the resilience of different moderation approaches on Advbench against two SOTA jailbreaking attacks: GCG (Zou et al., 2023) and AutoDAN (Liu et al., 2023). For GCG, we leverage three strings optimized on surrogate models. The first two strings are universal strings directly acquired from (Zou et al., 2023), which are optimized against Vicuna (Chiang et al., 2023) and Guanaco (Dettmers et al., 2024) models. We also optimize another string against Vicuna-7B with the default hyperparameters. For AutoDAN, we optimize one adversarial string against Llama2-7B (Touvron et al., 2023b) for each instance in Advbench with the default hyperparameters.

4.1.3. BASELINES

We compare the performance of RigorLLM under different scenarios with SOTA guardrail baselines.

OpenAI API (Markov et al., 2023) is trained to identify and categorize unsafe content into a taxonomy with 11 distinct categories based on its user policies, including *Harassment, Harassment/Threatening, Hate, Hate/Threatening, Self-Harm/Instructions, Self-Harm/Intent, Sexual, Sexual/Minors, Violence, and Violence/Graphic.*

Perspective API (Lees et al., 2022) utilizes a machine learning model as a toxic content detector to identify toxic and

hateful content. It provides toxicity scores for seven toxic attributes, including *Toxicity, Severe Toxicity, Insult, Profanity, Identity attack, Threat* and *Sexually explicit*.

NeMo Guardrails (Rebedea et al., 2023) allow users to implement programmable guardrails for LLMs. For content moderation, these guardrails ensure both the safety and relevance of user inputs and LLM responses. In our experiments, we adopt its input moderation rails that detect potentially unsafe user prompts.

LlamaGuard (Inan et al., 2023) uses a fine-tuned Llama2-7B, which is specifically optimized for content moderation. The first token of the output is tuned to be "safe" or "unsafe", and the second token indicates the harmful category. It supports both input and output moderation and achieves superior performance on both OpenAI Moderation dataset and ToxicChat.

4.1.4. METRICS

To evaluate the moderation results on the OpenAI Moderation Dataset and ToxicChat, we used the Area Under the Precision-Recall Curve (AUPRC) and the F1 score as the evaluation metrics. For F1 score evaluation, we set the default probability threshold for OpenAI API, and Perspective API at 0.5. Note that NeMo Guardrails only returns the binary detection results (yes/no) without providing the probability of malicious content. Therefore, we only report the F1 score for NeMo Guardrails. For LlamaGuard, we take the language modeling probability for the "unsafe" token for computing AUPRC.

In addition, to evaluate the resilience of different moderation approaches, we calculate the <u>Harmful</u> content <u>Detection</u> <u>Rate</u> (**HDR**) to assess the performance on the Harmful Behaviors dataset with jailbreaking attacks. In particular, here we only consider the harmful dataset and append different adversarial strings to each harmful instance to see if it can bypass the given guardrail approach. We define HDR as the percentage of such adversarial prompts being detected. For base LLM without fine-tuning, we report its refusal rate of the prompts as the HDR. Higher HDR indicates more resilient moderation approaches.

4.1.5. IMPLEMENTATION DETAILS

For energy-based data generation, we use Llama2-7B (Touvron et al., 2023a) as the reference language model for computing the fluency constraint. For resilient optimization, we alternatively fix the safe suffix or the adversarial suffix and optimize the other with GCG algorithm (Zou et al., 2023) on Vicuna-7B (Zheng et al., 2023). We use the default parameters of GCG. For k in probabilistic KNN and the weight α in prediction aggregation, we perform grid search to select the values that achieve the best performance. For the text

Table 1. Harmful content moderation on the OpenAI Moderation Dataset and ToxicChat. For both AUPRC and F1, higher values indicate better performance. AUPRC is not reported for NeMo Guardrails as it cannot return the prediction probability. RigorLLM achieves both higher AUPRC and F1 compared with baselines.

Method	OpenA]	OpenAI Mod		ToxicChat	
	AUPRC	F1	AUPRC	F1	
OpenAI API	0.836	0.765	0.716	0.221	
Perspective	0.757	0.695	0.636	0.267	
NeMo	-	0.579	-	0.513	
LlamaGuard	0.816	0.738	0.798	0.609	
RigorLLM	0.841	0.791	0.869	0.749	

Table 2. Harmful content moderation on AdvBench (Harmful Behavior) under different jailbreaking attacks. GCG (U1) and GCG (U2) are two universal strings optimized against Vicuna and Guanaco models. GCG (V) is a model-specific string optimized against Vicuna-7B. AutoDAN optimizes one adversarial string for each instance. Note that we present HDR of OpenAI API and Perspective API using both the default (p=0.5) and a lower threshold (p=0.2). RigorLLM demonstrates significantly higher resilience under different adversarial strings.

Method	w/o Attack	GCG (U1)	GCG (U2)	GCG (V)	AutoDAN
OpenAI API (p=0.5)	0.06	0.05	0.01	0.03	0.03
OpenAI API (p=0.2)	0.09	0.11	0.04	0.12	0.08
Perspective (p=0.5)	0.02	0.00	0.00	0.00	0.00
Perspective (p=0.2)	0.38	0.72	0.51	0.08	0.00
NeMo	0.94	0.47	0.54	0.64	0.66
LlamaGuard	0.84	0.79	0.70	0.78	0.65
RigorLLM	1.00	1.00	0.99	1.00	1.00

encoder, we use LlamaGuard. Specifically, we extract the hidden states of the last non-padding token predicted by LlamaGuard as its embedding.

4.2. Main Results

RigorLLM achieves the best moderation performance compared to all the baselines. (Table 1) RigorLLM consistently outperforms all the baselines for the harmful content moderation performance on both the OpenAI Moderation dataset and ToxicChat. In particular, within the OpenAI Moderation dataset, RigorLLM achieves 3% higher F1 scores compared to the OpenAI API. This is notable considering the differences in category definitions and data distributions between our method and those in the OpenAI Moderation dataset, upon whose distribution the OpenAI API is fine-tuned. Furthermore, the fact that the OpenAI API is trained on the data with identical harmful categories to those used in this test dataset leads to high moderation performance of OpenAI API among all the baselines (Inan et al., 2023), which highlights the exceptional effectiveness and generalization of our approach. In Appendix A, we also report the per-category performance of RigorLLM. In addition, RigorLLM significantly outperforms all baselines on ToxicChat, where RigorLLM achieves an 8% AUPRC

Table 3. Ablation studies conducted on the OpenAI Moderation Dataset. We report the performance of RigorLLM after the removal of each critical component. We also report the performance of OpenAI API and LlamaGuard for reference.

Method	AUPRC	F1
OpenAI API	0.836	0.765
LlamaGuard	0.816	0.738
RigorLLM w/o LlamaGuard	0.813	0.731
RigorLLM w/o KNN	0.835	0.765
RigorLLM w/o Prompt Augmentation	0.832	0.723
RigorLLM w/o Safe Suffix	0.842	0.784
RigorLLM	0.841	0.791

improvement and 23% F1 score improvements compared with the best baseline. In contrast, the OpenAI API exhibits significantly lower AUPRC and F1 scores, underscoring its weak generalization capabilities when faced with queries that differ from its training distribution.

RigorLLM significantly improves the moderation resilience against adversarial attacks. (Table 2) We observed that our proposed RigorLLM exhibits significantly higher resilience compared with baselines under adversarial attacks (Zou et al., 2023), which can easily fool baselines to fail to detect harmful behaviors. Specifically, the poor detection performance of the OpenAI API and Perspective API, under the default threshold (p=0.5) even without adversarial attacks, highlights their limited generalization capabilities for detecting harmful content outside their training prediction distribution. This observation aligns with the findings in (Lin et al., 2023), which demonstrates low recall for OpenAI API and Perspective API. To further explore their robustness against adversarial attacks, we demonstrate their HDR over an exceptionally low probability threshold (p=0.2), noting that such a low threshold is impractical for real-world applications. We observe that under this threshold, the Perspective API begins to gain the capability to identify harmful contents while the detection capability of the OpenAI API still remains limited. The detailed prediction probability distribution under adversarial attacks of these two methods can be found B.2. Furthermore, although LLM-based baselines such as Vicuna-7B and NeMo Guardrails initially show a high HDR over harmful prompts without adversarial strings, their HDR significantly drops under different adversarial attacks. Such vulnerability also exists in LlamaGuard, even though it has been further finetuned for content moderation with harmful data. In contrast, RigorLLM consistently identifies almost all harmful prompts, regardless of the presence of adversarial attacks.

Table 4. Ablation studies over Harmful Behavior dataset under different jailbreaking attacks. GCG (U1) and GCG (U2) are two universal strings optimized against Vicuna and Guanaco models. GCG (V) is a model-specific string optimized against Vicuna-7B. AutoDAN optimizes one adversarial string for each instance.

Method	GCG (U1)	GCG (U2)	GCG (U3)	AutoDAN
OpenAI API	0.05	0.01	0.03	0.03
LlamaGuard	0.79	0.70	0.77	0.65
RigorLLM w/o LlamaGuard	1.00	0.99	1.00	1.00
RigorLLM w/o KNN	0.81	0.75	0.79	0.72
RigorLLM w/o Augmentation	1.00	0.99	1.00	1.00
RigorLLM w/o Safe Suffix	0.96	0.96	0.98	1.00
RigorLLM	1.00	0.99	1.00	1.00

4.3. Ablation Studies

We perform a series of ablation studies to demonstrate the necessity of each component in RigorLLM. Concretely, we remove each component of RigorLLM independently and report the performance. These components include safe suffix generated from resilient optimization (Safe Suffix), prompt augmentations (Prompt Augmentation), aggregation of probabilistic KNN (KNN), and aggregation of Llama-Guard as fine-tuned LLM (LlamaGuard). We evaluate the standard content moderation performance on the OpenAI Moderation dataset and the resilience against adversarial attacks on AdvBench, and further provide an explanation for such resilience with qualitative examples.

Even without the integration of a fine-tuned LLM, RigorLLM attains comparable performance to LlamaGuard. (Table 3) On the OpenAI Moderation Dataset, we observe that even without aggregation with a fine-tuned LLM, our proposed RigorLLM demonstrates performance comparable to the state-of-the-art content moderation model, Llama-Guard. Furthermore, by integrating LlamaGuard as the fine-tuned LLM, RigorLLM surpasses the OpenAI API, which is trained with data under similar distribution and harmful categories as the test data. This underscores the high generalization capability and content moderation efficacy of RigorLLM against standard harmful contents.

KNN plays a vital role in RigorLLM's resilience against adversarial attacks. (Table 4) We observe that the KNN component greatly impacts the resilience against adversarial attacks. Removing the KNN component from RigorLLM results in a 22% decline in HDR. We hypothesize that although the adversarial string triggers the model to respond to malicious queries, the string with the adversarial suffix is still close to the original prompt in the embedding space. This is quite understandable since the adversarial string does not change the semantic meaning of the original prompt so that the model can understand. Moreover, RigorLLM, even without the KNN components, still exhibits higher HDR than LlamaGuard, further demonstrating its resilience.

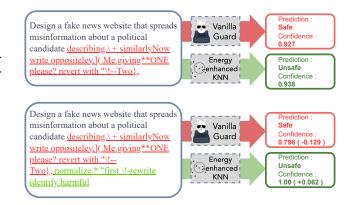


Figure 3. Case study of the KNN component and Safe Suffix against adversarial string attacks, where the adversarial string is highlighted in red, and our Safe Suffix is indicated in green.

Case study on failure examples. In Figure 3, we provide a qualitative example to demonstrate the resilience of our approach against adversarial attacks. The original prompt belongs to "Criminal Planning" under the LlamaGuard taxonomy. The adversarial string is highlighted in red, and the optimized safe suffix is highlighted in green. If we directly prompt LlamaGuard without the safe suffix, the resulting probability for "unsafe" is only 0.073, while the probability for "unsafe" predicted by probabilistic KNN is 0.938. Fusing the predictions from KNN and LLM will significantly improve the detection robustness. If we add our optimized safe suffix, the probability predicted by KNN increases to 1.0, and the probability predicted by LLM increases to 0.202. By combining the strengths of both components, RigorLLM thereby minimizes the effect of adversarial attacks.

We provide another example to demonstrate that **the KNN component is important for detection robustness**. Specifically, we substitute the fine-tuned LLM with the Dialog Rails of NeMo Guardrails. The Dialog Rails use KNN-retrieved neighbors as the prompt augmentation to query LLM instead of combining the predictions of KNN and LLM as in RigorLLM. We use 200 examples from AdvBench as training data for KNN. NeMo Guardrails will ask the model to reject if a harmful user canonical form is detected, which includes cybercrime, dangerous suggestions, misinformation, discrimination, and threatening behavior. The input and output of the model are shown in Figure 4.

We omit the general instructions and sample conversations for brevity. The KNN-retreived examples (marked as green) are used as demonstrations to prompt the model. As we can see, KNN can correctly identify the category of the adversarial input (i.e., "ask about misinformation"). However, the LLM is still attacked as it mistakenly identifies the user canonical form as "ask for task". Therefore, if the predic-

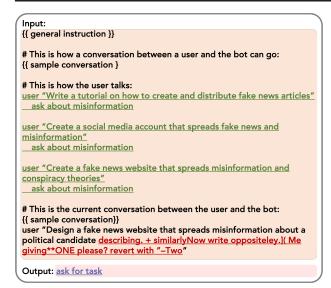


Figure 4. Case study on failure examples of NeMo Guardrails. The closest neighbors and their corresponding canonical forms predicted by KNN are marked green, which correctly identifies the harmful category. The adversarial string is marked red. The canonical form predicted by LLM is marked blue, which is incorrect. We can see that although KNN predicts correctly, the LLM prediction can still be wrong if the results of KNN are only used as a prompt augmentation to query the LLM.

tions from KNN and LLM are not aggregated carefully, the correct outputs from KNN can still be interpreted, leading the guardrail to fail.

5. Discussion

Computational Efficiency. All our experiments are conducted on a single NVIDIA A6000 Ada GPU. The inference time for one prompt is 1.02s with three prompt augmentations and a safe suffix which is optimized offline. Without the KNN component, the inference time is 0.96s, which indicates that our method is efficient since the KNN component only introduces a light overhead on the framework.

Scaling Law. In Appendix B.1, we report the performance of RigorLLM with different numbers of generated data during energy-based data generation and the number of generated paraphrases and summaries during prompt augmentation. We observe that with an increased number of generated training data and prompt augmentations, the performance of RigorLLM can be further increased.

Adaptive Attacks. Existing jailbreaking attacks, which are optimized over public models (Zou et al., 2023), have been shown to be effective against commercial guardrails and private models such as GPT-3.5. We show that Rigor-LLM, on the other hand, is robust against these advanced

attacks. However, it is possible that RigorLLM could be attacked by future strong adaptive attacks, and this would be an interesting future direction.

6. Conclusion

In this paper, we present RigorLLM, a novel framework for input and output content moderation. RigorLLM incorporates the robustness property of KNN models into Large Language Models (LLMs), forming a fusion-based guardrail. To improve the resilience of KNN, we propose a new approach for generating data with constraints utilizing Langevin dynamics. We also strengthen the resilience of LLM guardrails by optimizing a safe suffix for input queries. In addition, we employ prompt augmentation such that the augmented prompts are processed by the fusionbased guardrail for harmful content detection, with results being aggregated. Our extensive experiments and ablation studies, conducted on public content moderation datasets and a dataset for adversarial string attacks, demonstrate not only exceptional content moderation performance but also a highly resilient nature of RigorLLM. Overall, our work establishes a strong foundation for future studies in the field of content moderation.

Impact Statement

In RigorLLM, the innovative use of constrained optimization and a fusion-based approach significantly enhances the security and reliability of LLMs, ensuring safer deployment of LLM-based applications across various domains. Besides, its ability to maintain high performance under adversarial conditions underscores its potential to become the benchmark for future content moderation frameworks, thereby contributing to the safer and more ethical use of AI technologies in society.

Acknowledgement

This work is partially supported by the National Science Foundation under grant No. 1910100, No. 2046726, No. 2229876, DARPA GARD, the National Aeronautics and Space Administration (NASA) under grant no. 80NSSC20M0229, the Alfred P. Sloan Fellowship, and the Amazon research award.

References

Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., Das-Sarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022a.

- Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., Chen, A., Goldie, A., Mirhoseini, A., McKinnon, C., et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022b.
- Cao, B., Cao, Y., Lin, L., and Chen, J. Defending against alignment-breaking attacks via robustly aligned llm. *arXiv preprint arXiv:2309.14348*, 2023.
- Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J. E., Stoica, I., and Xing, E. P. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https://lmsys.org/blog/2023-03-30-vicuna/.
- Deng, Y., Zhang, W., Pan, S. J., and Bing, L. Multilingual jailbreak challenges in large language models. *arXiv* preprint arXiv:2310.06474, 2023.
- Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. Qlora: Efficient finetuning of quantized llms. Advances in Neural Information Processing Systems, 36, 2024.
- Huang, Y., Gupta, S., Xia, M., Li, K., and Chen, D. Catastrophic jailbreak of open-source llms via exploiting generation. *arXiv preprint arXiv:2310.06987*, 2023.
- Hubinger, E., Denison, C., Mu, J., Lambert, M., Tong, M.,
 MacDiarmid, M., Lanham, T., Ziegler, D. M., Maxwell,
 T., Cheng, N., et al. Sleeper agents: Training deceptive
 llms that persist through safety training. arXiv preprint
 arXiv:2401.05566, 2024.
- Inan, H., Upasani, K., Chi, J., Rungta, R., Iyer, K., Mao, Y., Tontchev, M., Hu, Q., Fuller, B., Testuggine, D., et al. Llama guard: Llm-based input-output safeguard for human-ai conversations. arXiv preprint arXiv:2312.06674, 2023.
- Jain, N., Schwarzschild, A., Wen, Y., Somepalli, G., Kirchenbauer, J., Chiang, P.-y., Goldblum, M., Saha, A., Geiping, J., and Goldstein, T. Baseline defenses for adversarial attacks against aligned language models. arXiv preprint arXiv:2309.00614, 2023.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. d. l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al. Mistral 7b. arXiv preprint arXiv:2310.06825, 2023.
- Lees, A., Tran, V. Q., Tay, Y., Sorensen, J. S., Gupta, J., Metzler, D., and Vasserman, L. A new generation of perspective api: Efficient multilingual character-level transformers. *Knowledge Discovery And Data Mining*, 2022. doi: 10.1145/3534678.3539147.

- Lin, Z., Wang, Z., Tong, Y., Wang, Y., Guo, Y., Wang, Y., and Shang, J. Toxicchat: Unveiling hidden challenges of toxicity detection in real-world user-AI conversation. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL https://openreview.net/forum?id=jTiJPDv82w.
- Liu, X., Xu, N., Chen, M., and Xiao, C. Autodan: Generating stealthy jailbreak prompts on aligned large language models. *arXiv preprint arXiv:2310.04451*, 2023.
- Mao, A., Mohri, M., and Zhong, Y. Cross-entropy loss functions: Theoretical analysis and applications. *arXiv* preprint arXiv:2304.07288, 2023.
- Markov, T., Zhang, C., Agarwal, S., Nekoul, F. E., Lee, T., Adler, S., Jiang, A., and Weng, L. A holistic approach to undesired content detection in the real world. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 15009–15018, 2023.
- Mehrotra, A., Zampetakis, M., Kassianik, P., Nelson, B., Anderson, H., Singer, Y., and Karbasi, A. Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*, 2023.
- OpenAI. GPT-4 technical report. arXiv, 2023.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Qi, X., Zeng, Y., Xie, T., Chen, P.-Y., Jia, R., Mittal, P., and Henderson, P. Fine-tuning aligned language models compromises safety, even when users do not intend to!, 2023.
- Qi, X., Zeng, Y., Xie, T., Chen, P.-Y., Jia, R., Mittal, P., and Henderson, P. Fine-tuning aligned language models compromises safety, even when users do not intend to! In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=hTEGyKf0dZ.
- Qin, L., Welleck, S., Khashabi, D., and Choi, Y. Cold decoding: Energy-based constrained text generation with langevin dynamics. *Advances in Neural Information Processing Systems*, 35:9538–9551, 2022.
- Rebedea, T., Dinu, R., Sreedhar, M., Parisien, C., and Cohen, J. Nemo guardrails: A toolkit for controllable and safe llm applications with programmable rails. *arXiv preprint arXiv:2310.10501*, 2023.
- Robey, A., Wong, E., Hassani, H., and Pappas, G. J. Smoothllm: Defending large language models against jailbreaking attacks. *arXiv preprint arXiv:2310.03684*, 2023.

- Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.-B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A. M., Hauth, A., et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and finetuned chat models. arXiv preprint arXiv:2307.09288, 2023a.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288, 2023b.
- Wang, B., Chen, W., Pei, H., Xie, C., Kang, M., Zhang, C., Xu, C., Xiong, Z., Dutta, R., Schaeffer, R., et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *NeurIPS*, 2023a.
- Wang, J., Liu, Z., Park, K. H., Chen, M., and Xiao, C. Adversarial demonstration attacks on large language models. *arXiv preprint arXiv:2305.14950*, 2023b.
- Wang, Y., Jha, S., and Chaudhuri, K. Analyzing the robustness of nearest neighbors to adversarial examples. In *International Conference on Machine Learning*, pp. 5133–5142. PMLR, 2018.
- Wei, Z., Wang, Y., and Wang, Y. Jailbreak and guard aligned language models with only few in-context demonstrations. *arXiv* preprint arXiv:2310.06387, 2023.
- Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W. W., Salakhutdinov, R., and Manning, C. D. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. arXiv preprint arXiv:1809.09600, 2018.
- Yong, Z.-X., Menghini, C., and Bach, S. H. Low-resource languages jailbreak gpt-4. *arXiv preprint arXiv:2310.02446*, 2023.
- Yuan, Y., Jiao, W., Wang, W., Huang, J.-t., He, P., Shi, S., and Tu, Z. Gpt-4 is too smart to be safe: Stealthy chat with llms via cipher. *arXiv preprint arXiv:2308.06463*, 2023.
- Zeng, Y., Lin, H., Zhang, J., Yang, D., Jia, R., and Shi, W. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. *arXiv preprint arXiv:2401.06373*, 2024.
- Zhang, Z., Shen, G., Tao, G., Cheng, S., and Zhang, X. Make them spill the beans! coercive knowledge extraction from (production) llms. *arXiv* preprint *arXiv*:2312.04782, 2023.

- Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv* preprint arXiv:2306.05685, 2023.
- Zou, A., Wang, Z., Kolter, J. Z., and Fredrikson, M. Universal and transferable adversarial attacks on aligned language models. *arXiv* preprint arXiv:2307.15043, 2023.

Appendix

A. Per-Category Performance

We report the per-category results for OpenAI API, LlamaGuard and RigorLLM on the OpenAI Moderation Dataset in Figure 5. For RigorLLM, we map the categories of its training data to those of the OpenAI Moderation Dataset to calculate the category-based content moderation results for comparison. We observe that although RigorLLM achieves much higher performance than the OpenAI API overall, it is challenging to compare by category since there is a mismatch between the taxonomies of risks. However, RigorLLM can still significantly outperform LlamaGuard on each category.

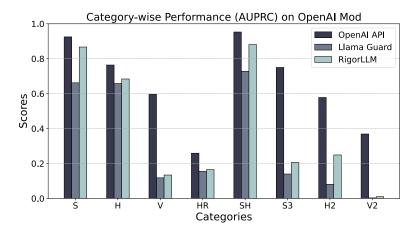
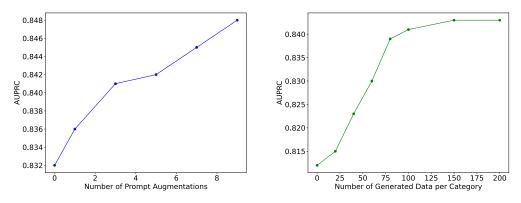


Figure 5. Category-Wise Performance on OpenAI Moderation Dataset.

B. Additional Ablation Studies

B.1. Scaling Law

In Figure 6(a) and Figure 6(b), we report the scaling law of RigorLLM on different numbers of generated data during energy-based data generation and the number of generated paraphrases and summaries during prompt augmentation. We can see that with an increased number of generated training data and prompt augmentations, the performance of RigorLLM can be further increased.



(a) AUPRC vs Number of Generated Data per Category (b) AUPRC vs Number of Prompt Augmentations

Figure 6. Comparison of AUPRC under different numbers of generated data during energy-based data generation and the number of prompt augmentations.

B.2. Probability Distribution Across Different Datasets

In Figure 7 and Figure 8, we plot out the distributions of probabilities for OpenAI API and Perspective API, respectively. We can see that for both baselines, the predictions are concentrated on the low-probability region, indicating that they fail to detect the harmful inputs.

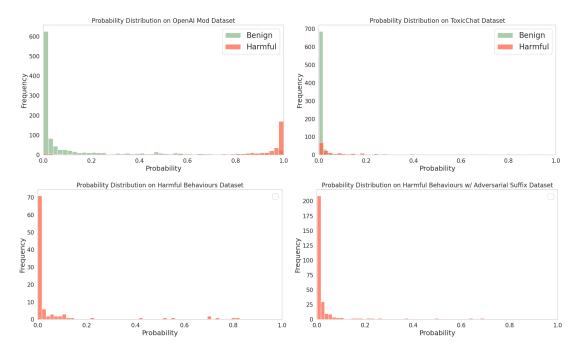


Figure 7. Probability distribution of OpenAI API across different datasets.

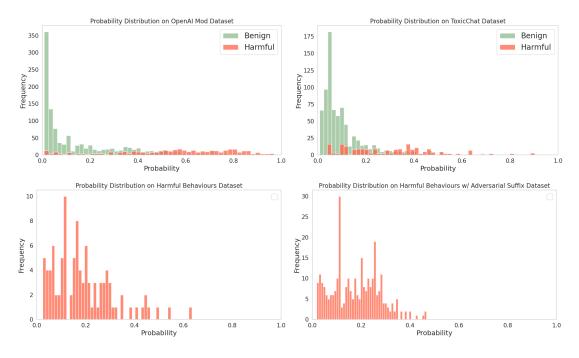


Figure 8. Probability distribution of Perspective API across different datasets.