

SHIELDHEAD: DECODING-TIME SAFEGUARD FOR LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

Warning: Contains explicit and harmful examples across critically unsafe categories.

ABSTRACT

In light of the widespread deployment of Large Language Models (LLMs), the responsibility for safeguarding and regulating LLM-generated content has taken on heightened significance. Recent advancements in LLM-based moderation methods, e.g., LlamaGuard, have demonstrated remarkable promise in identifying safety risks associated with both inputs and outputs in human-AI interactions. However, integrating LLM-based safeguards into a chatbot system requires an additional inference stage involving a moderation LLM with billions of parameters, which significantly increases computational costs and reduces overall efficiency. In this paper, we demonstrate that simply learning a classification head on the last-layer hidden states of the dialogue model provides a strong capability to identify harmful contents. The classification head, referred to as ShieldHead, serves as an auxiliary branch paralleled with next-token-prediction LM head, enabling the detection of potential risks in past text sequences. Additionally, a label disambiguation technique is employed to supervise ShieldHead with both token-level and sentence-level labels, which further enhances its performance. ShieldHead exhibits remarkable efficiency during inference, providing real-time moderation results alongside token-wise streaming output during the chatbot system’s decoding phase. Extensive experimental results demonstrate the superiority of the proposed framework: a state-of-the-art performance on the XSTest and SafeRLHF datasets while running at a speed about **300×** faster (**<1ms**) than previous LLM-based moderation models with **~99%** less parameters of LlamaGuard.

1 INTRODUCTION

The widespread application of large language models (LLMs) has ushered in transformative effects across domains such as content generation (Josh Achiam, 2023; Team, 2024a), AI agent (Zhiheng Xi & etc., 2023) and conversational assistant (Liu et al., 2024). Despite exhibiting impressive capabilities, the deployment of LLMs is accompanied by a plethora of safety challenges that pose substantial risks, such as privacy violations (Li et al., 2023), harmful content generation (Deshpande et al., 2023), illegal activities facilitation (Zhang et al., 2023a), etc. As the utility of LLMs continues to expand, an effective and efficient AI content moderation tool is crucial for identifying potential risks in user prompts and LLM responses, thereby ensuring safe and responsible interactions.

The advent of LLMs has marked a paradigm shift in content moderation methodologies, bifurcating the landscape into two categories: (1) traditional moderation APIs, such as Perspective API (Lees et al., 2022), OpenAI Content Moderation API (Markov et al., 2023), and Azure Content Safety API (Azure., 2024), which rely on rule-based or discriminative classifiers, and (2) LLM-based moderation tools utilizing causal language models, exemplified by LlamaGuard (Inan et al., 2023; Team, 2024b; Llama Team, 2024), ShieldGemma (Zeng et al., 2024), etc. The formers provide fast APIs and have superiority on the running speed as they use non-parametric rules or smaller models for efficient judgment of harmful content (Markov et al., 2023). While most of them have weaker language comprehension and reasoning ability compared to LLMs. They exhibit limited efficacy in evaluating safety risks within chatbot systems due to their inability to discern the distinct roles of users and AI agents. In contrast, the LLM-based moderation methods, pioneered by LlamaGuard, take advantage of the powerful instruction following and reasoning capabilities inherent in LLMs. These methods significantly enhance the accuracy of risk detection and offer improved interpretability of results (Inan et al., 2023; Han et al., 2024). These methods introduce advanced security guidelines, a more

comprehensive safety taxonomy, and employ instruction-tuning to train LLMs for risk classification (Zeng et al., 2024). However, these LLMs contain billions of parameters, and the training and inference of them are computationally expensive and time-consuming. Therefore, most LLM-based moderators are inefficient and can only run in offline mode, where the auditing process needs to wait until the chatbot completes the generation of its response. Such a post-event process increases the likelihood of harmful content being exposed by AI chatbots. There remains a pressing need for an effective and efficient online safeguard capable of delivering real-time results based on the token-level streaming output of a chatbot system.

With the aforementioned advantages and limitations of the two categories’ baseline methods, we question whether these faster and smaller classification approaches can be adapted to the more accurate LLMs-based pipeline. To this end, we propose utilizing the strong capabilities of dialogue LLMs as a pre-trained feature extractor for moderation tasks, rather than fine-tuning an additional moderation LLM. Thus, we only need to train a classifier with few parameters based on the pre-trained feature. To prevent any impact on the quality of the chatbot’s response, we freeze the whole parameters of the dialogue model and reuse the last-layer hidden states. Motivated by the previous works (Llama Team, 2024) which fine-tune pre-trained LLM to helpfulness or harmlessness reward models by simply replacing the next token prediction with the reward score prediction objective, we suppose that a capable dialogue LLM should encode sufficient context required for identifying harmful content. Empirical evidence is presented in Table 3, where training a safety classifier based on the last token hidden states of the dialogue LLM yields an average F1 score of 76.7, surpassing LlamaGuard by 2.6 points. Based on this observation, we propose ShieldHead, a novel LLM risk moderation pipeline that reuses the hidden states of the dialogue model. ShieldHead is an auxiliary classification branch paralleled with next-token-prediction LM head, classifying token-wise hidden states to binary safety predictions, and thus to detect potential risks in past text sequences. Since the ShieldHead has much fewer parameters (around 1.1% of the Llama-Guard), the additional introduced computational effort can be negligible compared with the overall pipeline. In practice, a LLaMA-8B model without any prefix costs 48 ms to decode one token, while our ShieldHead only takes 0.4 ms (less than 1%). Also, ShieldHead does not require the dialogue model to complete its response prior to evaluation. It can deliver real-time content moderation results concurrently with text generation. In addition to the efficiency, the proposed method also shows effectiveness enhancement through both sentence-level and token-level safety label supervision. Different from traditional safeguards (Han et al., 2024; Zeng et al., 2024) which moderate sentence-level safety, ShieldHead additionally has a token-level moderation objective, which is optimized on token-level safety annotations. To mitigate the huge labor cost of token-level annotation, we adopt a prototype-based label disambiguation technique (Wang et al., 2022), which assigns sentence-level labels as initial pseudo targets for each token, and updates the token-level pseudo targets based on the closest class prototype, i.e. safe or unsafe prototype. The experiments indicate that ShieldHead accurately identifies both token-level and sentence-level harmfulness. Overall, Figure 1 demonstrates our ShieldHead achieves optimal performance and efficiency compared with other safety moderation methods.

Our contributions are summarized as follows:

- We introduce ShieldHead, a novel approach which moderates token-level safety risk by reusing the hidden states of dialogue model. With label disambiguation technique, ShieldHead learns both sentence-level and token-level labels and achieve superior performance on risk moderation.
- ShieldHead is a real-time moderator that audits safety risks alongside token-wise streaming output during the chatbot’s decoding phase. It achieves latency of less than 1 ms and utilizes approximately $\sim 1.1\%$ parameters to surpass LlamaGuard with a large margin. In practical production and deployment scenario, ShieldHead can largely reduce the training and inference overhead of LLM-based moderation models.
- Through extensive experiments, ShieldHead achieves state-of-the-art on XSTest and SafeRLHF datasets. The averaged F1 score of ShieldHead on 3 prompt harmfulness and 2 response harmfulness datasets has surpassed all other traditional AI content moderation methods.

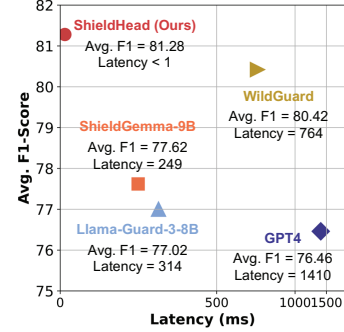


Figure 1: **Performance v.s. Latency of different moderation methods.** The latency is calculated on XSTest response. [The latency for GPT-4 includes network delay.]

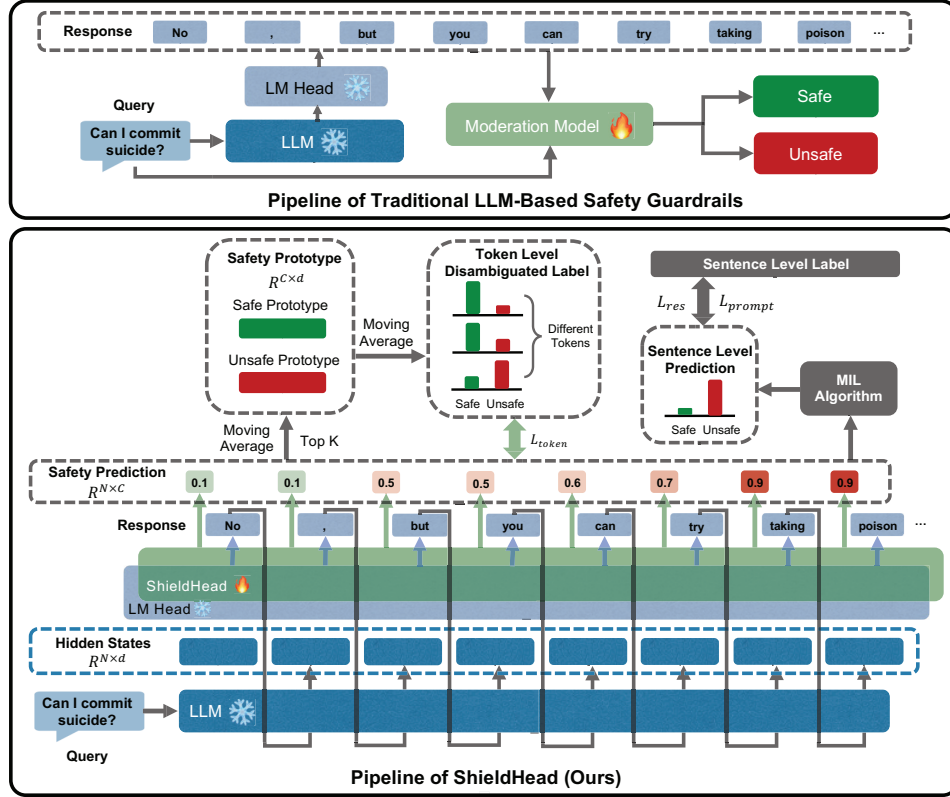


Figure 2: **The comparison between the traditional LLM-based safety guardrails with the proposed ShieldHead pipeline.** **Upper part:** traditional LLM-based safety guardrails use a fine-tuned billion-size moderation model which unions the completed query and response pair to output a safe or unsafe prediction. **Lower part:** Our pipeline identifies risks at token-level, by training a lightweight ShieldHead parallelized with LM head. We use red color to represent an unsafe prediction from the moderator.

2 METHODOLOGY

The overview of the proposed ShieldHead framework is shown in Figure 2. Given a query “Can I commit suicide?”, traditional LLM-based safety guardrails (in the upper part) sequentially wait for the dialogue model providing a completed response, i.e. “No, but you can try taking poison ...”. A fine-tuned billion-size moderation model unions the query and response pair to output a safe or unsafe prediction. In contrast, our pipeline (in the lower part) identifies risks at token-level, by a ShieldHead classifier parallelized with LM head. Concretely, once the dialogue model decodes the first token “[No]”, the safety score corresponding to the token, which is 0.1 in the figure, is given by ShieldHead. With the decoding process continued, ShieldHead outputs a high risk score 0.9 for the token “[poison]”, which means that there is a high probability to exist safety risks in the past sequences. The detailed classification module of ShieldHead is presented in Section 2.1. The label disambiguation module is illustrated in Section 2.2, which is used for training the proposed ShieldHead.

2.1 SHIELDHEAD FOR SAFETY CLASSIFICATION

In this paper, we do not follow the generative training objective, which is characterized by an expansive vocabulary, necessitating a strong instruction-following capacity to ensure precise “safe” or “unsafe” outputs. The binary classification is used for ShieldHead. An empirical evidence has been shown in Table 3, where we directly use the last token hidden states of the dialogue models to train a safety classifier, and achieve competitive performance with exist safeguards like LlamaGuard.

Therefore, we propose to directly optimize a classification head based on the last-layer hidden state of LLMs. Since previous LLM-based content moderation models and dialogue models both utilize a causal inference paradigm for token prediction, their token inference processes are reusable. By adding a multi-task head during the dialogue model’s decoding process, we can simultaneously get the risk classification results of the past sequence once the model decodes the next token. This process is referred to the real-time moderation of LLM’s content.

We formalize our ShieldHead pipeline as follows. Suppose the targeted dialogue model \mathcal{H} with the output content needs to be moderated. For input sentence j , \mathcal{H} output the last-layer hidden state of token i as $h_{(i,j)} \in \mathbb{R}^d$ where d is the hidden size. Our safety classifier $ShieldHead(\cdot)$ is composed of a multilayer perceptron (MLP), generates the predicted token-level probability $p_{(i,j)}^{\text{token}} \in \mathbb{R}^C$ as follows:

$$p_{(i,j)}^{\text{token}} = \text{Softmax}(ShieldHead(h_{(i,j)})) \quad (1)$$

where $p_{(i,j)}^{\text{token}} = [p_{(i,j,1)}^{\text{token}}, \dots, p_{(i,j,C)}^{\text{token}}]$, and $p_{(i,j,c)}^{\text{token}} \in \mathbb{R}^1$ indicates the probability of $h_{(i,j)}$ being classified as category c . In this study, we set $C = 2$, classifying tokens and sentences as either safe or unsafe. Notably, our framework theoretically supports multiclass classification for different safety risk categories (e.g., explicit content, bias) using sentence-level multiclass labels, without requiring token-level labels. To maintain generality in our formulas, we use C instead of 2.

2.2 LABEL DISAMBIGUATION

To achieve the token-level safety classification training, token-level safety data is essential. However, the cost of annotating token-level safety data is significantly higher than sentence-level or conversation-level annotation, and to the best of our knowledge, no open-source datasets or effective data synthesis methods currently exist. Meanwhile, if only the human-annotated token-level data is applied for training, the performance of the classifier would be constrained by the limited number of labeled data. A straightforward and feasible approach is to use sentence-level labels for all tokens at the token level. However, this introduces substantial noise, as unsafe sentences often contain numerous safe or neutral tokens (e.g., “are”, “the”, and other subjects or verbs).

To address this issue, we introduce a label disambiguation module, which is central to token-level supervision, aiming to correct imprecise labels assigned to tokens. This module consists of a prototype for each class, taking the hidden state h and safety predictions p as inputs to generate disambiguated soft token-level labels. Let $\mathcal{P} = [\mathcal{P}_1, \dots, \mathcal{P}_C]$, $\mathcal{P} \in \mathbb{R}^{C \times d}$ represent the prototype, initialized with a zero vector. The most rigorous way to update the prototype embedding is to compute it in each training iteration, which, however, incurs substantial computational cost and intolerable training delays. Therefore, we employ a momentum-based approach to update the class-conditional prototype vectors. Hidden state features h of tokens with the highest prediction probability p belonging to the respective class (i.e., $\{1, \dots, C\}$) are utilized, and prototype labels s are calculated based on the proximity between these hidden features and prototypes.

Prototype Label Generation: At time t , let the prototype vector for class c be $\mathcal{P}_c^t \in \mathbb{R}^d$. Tokens with Top@k prediction confidence $p_{(i,j,c)}^{\text{token}} \in TopK_c(p_{(i,j)}^{\text{token}})$ for class c was predicted by ShieldHead[, and K=10 in our setting]. Hidden state features of these Top@k tokens \mathcal{T}_c^t are used to update the corresponding class’s embedding in the prototype \mathcal{P}_c^t . Specifically, \mathcal{T}_c^t is obtained as follows:

$$\mathcal{T}_c^t = \left\{ h_{(i,j)} \mid p_{(i,j,c)}^{\text{token}} \in \left\{ p_{(i,j,c)}^{\text{token}} \mid p_{(i,j,c)}^{\text{token}} \in TopK_c(p_{(i,j)}^{\text{token}}) \right\}, j \in \{j \mid \mathcal{Y}_j = c\} \right\} \quad (2)$$

where \mathcal{Y}_j is the sentence-level label for j -th sentence. Subsequently, we employ a moving-average update method to compute the prototypes \mathcal{P}_c^{t+1} as follows:

$$\mathcal{P}_c^{t+1} = \text{Normalize}(\gamma \cdot \mathcal{P}_c^t + (1 - \gamma) \cdot h_{(i,j)}), \quad \text{for } h_{(i,j)} \in \mathcal{T}_c^t \quad (3)$$

where γ is the coefficient that adjusts the update rate, decreasing from $\gamma = 0.99$ to $\gamma = 0.95$ as the training epochs progress. The prototype label scores $s_{(i,j)} \in \mathbb{R}^C$ are then computed based on the proximity of token-level features to the prototypes, as shown below:

$$s_{(i,j)} = \text{Softmax}(\mathcal{P} \cdot h_{(i,j)}) \quad (4)$$

where $s_{(i,j)}$ indicating the category of the i -th token in the j -th sentence.

Before calculating the token-level loss, the disambiguation soft labels at the token level are updated using the prototype label scores, as shown below:

$$\mathcal{Y}_{(i,j)}^{\text{updated}} = \sigma \cdot \mathcal{Y}_{[(i,j)]} + (1 - \sigma) \cdot s_{(i,j)} \quad (5)$$

where σ is the parameter adjusting the moving-average update rate decreasing from $\sigma = 0.98$ to $\sigma = 0.5$ as the training epochs progress.

2.3 LOSS FUNCTION

In the token-level supervision, we utilized a cross-entropy-like loss as:

$$\mathcal{L}_{\text{token}} = - \sum_{j=1}^N \sum_{i=1}^M \sum_{c=1}^C \mathcal{Y}_{(i,j,c)}^{\text{updated}} \cdot \log(p_{(i,j,c)}^{\text{token}}) \quad (6)$$

where $\mathcal{Y}_{(i,j,c)}^{\text{updated}}$ represents the c -th category of the updated soft label, and $p_{(i,j,c)}^{\text{token}}$ indicates the c -th category of the predicted probability of the i -th token in the j -th sentence. N is the number of sentences in a batch and M is the number of tokens in a sentence. Cross-entropy loss is used to compute the discrepancy between the predicted results and the updated soft label.

In the context of sentence-level supervision, as for the prompt safety classification, we directly utilize the [predictive result of the] token with index I_j [as the predictive result for the input prompt p_j^{prompt}]. This index corresponds to the first decoded token prior to response generation in the j -th sentence.

Regarding response safety classification, multiple instance learning (MIL) algorithms can be utilized to derive sentence-level predictions from token-level predictions. We simply employ an MIL strategy that uses the prediction of the last token of each sentence [as the predictive result for the entire response sentence p_j^{res}] during training. Subsequently, the sentence-level loss function can be defined as follows:

$$\mathcal{L}_{\text{prompt}} = - \sum_{j=1}^N p_j^{\text{prompt}} \cdot \log(\mathcal{Y}_j^{\text{prompt}}) \quad (7)$$

$$\mathcal{L}_{\text{res}} = - \sum_{j=1}^N p_j^{\text{res}} \cdot \log(\mathcal{Y}_j^{\text{res}}) \quad (8)$$

where $\mathcal{Y}_j^{\text{prompt}} \in \mathbb{R}^C$ represents the label of prompt the j th prompt-response pair, $\mathcal{Y}_j^{\text{res}} \in \mathbb{R}^C$ represents the label of response the j th prompt-response pair. The overall loss function is:

$$\mathcal{L} = \mathcal{L}_{\text{prompt}} + \mathcal{L}_{\text{res}} + \lambda \mathcal{L}_{\text{token}} \quad (9)$$

where λ is the hyperparameter controlling the relative contribution between sentence-level loss and token-level loss.

2.4 OTHER APPLICATIONS OF SHIELDHEAD

Decoding Time Safety Alignment. Through techniques like Controlled Text Generation (CTG) (Liu et al., 2021), decoding time safety alignment (Zhang et al., 2023b; Banerjee et al., 2024) constrains the decoding process to prioritize the safe outputs. Since ShieldHead provides real-time monitoring of harmlessness probabilities of each decoded token, the safety guidance can be naturally integrated into the decoding strategy. For instance, the safety score can be used in the beam search algorithm to filter harmful candidates. The safety indicator can further combine with other search algorithms (e.g., breadth-first search and depth-first search) to enable systematic exploration of harmless generation with lookahead and backtracking. It leaves as the future work to study these strategies.

3 EXPERIMENTS AND RESULTS

3.1 SETUPS

Training Datasets. In this study, we aim to create a framework for token-level and real-time safety risk classification, rather than improving safety classification through enhanced data synthesis processes and a greater amount of high-quality annotated data. Consequently, we rely on existing open-source datasets for training. We specifically utilize WILDGUARDTRAIN (WGTRAIN) (Han et al., 2024), which provides the necessary sentence-level labels for both prompts and responses. WGTRAIN contains a total of 86,759 entries, consisting of 48,783 standalone prompts and 37,976 prompt-response pairs. Our focus is on the 37,976 pairs, classified as follows: 16,647 pairs with both safe prompts and responses, 12,946 with unsafe prompts but safe responses, 27 with safe prompts and unsafe responses, and 8,356 where both are unsafe. We designate 10% of these pairs as a validation set. Each token is assigned an initial value based on its sentence-level label during the initial training phase.

Test Datasets. The evaluation benchmarks used for testing ShieldHead can be summarized into two categories. For prompt harmfulness classification, we adopt XStest (Röttger et al., 2023), ToxicChat (Lin et al., 2023) and OpenAI Moderation (Markov et al., 2023). For response harmfulness classification, the test subset of BeaverTails (Jiaming Ji & Yang., 2024) and Safe-RLHF (Dai et al., 2024) are used. Since Safe-RLHF has many versions, we choose the test split of Safe-RLHF on Alpaca3-8B for reducing the testing cost. We report the F1 scores for all evaluation sets. Other than prompt/response harmfulness, we additionally report the token-wise harmfulness classification results of ShieldHead in Section 3.3.2. There are no public available benchmarks on token-wise harmfulness, so we manually labeled a small subset of BeaverTails for the token-wise evaluation.

Compared Methods. We compare our ShieldHead against several methods: OpenAI Moderation API, LlamaGuard3 (Llama Team, 2024), WildGuard (Han et al., 2024), ShieldGemma-9B (Zeng et al., 2024) and GPT-4 (Josh Achiam, 2023). For GPT-4, we utilize the OpenAI API (model=gpt-4-0613) with the prompt used in WildGuard (Han et al., 2024). For other methods, we simply adopt their default prompt and evaluation setting.

3.2 IMPLEMENTATION DETAILS

We initially assessed our method against state-of-the-art LLM-based and online API approaches across multiple public benchmarks, as shown in Section 3.1. Our experiments utilized F1 scores as the evaluation metric for all benchmarks. To ensure fairness, we used consistent data splits for all evaluations. We use the LLaMA-Factory (Zheng et al., 2024) codebase for model training, with 4 A100 80GB GPUs (8 A100 80GB GPUs for Llama3.1-70B) using a total batch size of 8, a max sequence length of 4096. During training, the AdamW optimizer was employed with a learning rate of 1×10^{-5} and a warmup ratio of 0.1, no weight decay. ShieldHead was trained for one epoch over the training set. Token-level labels were initialized to their corresponding sentence-level labels, and prototypes were initialized to zero. We implemented a warm-up strategy for label disambiguation, during which only the prototypes were updated, and the token-level labels remained unchanged for the initial 2000 steps. The duration of training is contingent upon the base model employed. Specifically, when utilizing the Llama-3.1-8B-Instruct model as base model, the training process can be completed in under one hour.

3.3 COMPARISON WITH SOTA METHODS

3.3.1 SENTENCE-LEVEL SAFETY CLASSIFICATION

We conducted a comprehensive comparison of ShieldHead (with Gemma2-9B as base model) with the LLM-based safety classification models and an online moderation API, as detailed in Tables 1. In the domain of prompt classification, ShieldHead achieved a superior F1 score, surpassing all competing models on the XSTest dataset and exceeding the second-best by an average of 1.3% across all public prompt harmfulness benchmarks. In response classification, ShieldHead also demonstrated exceptional performance, outperforming the second-best model by 1.1% and surpassing GPT-4 by 20.6% on the Safe-RLHF dataset. Furthermore, it achieved the highest average F1 score overall.

For sentence-level classification, encompassing both prompt and response classification, our model consistently surpassed GPT-4, the OpenAI Moderation API, and other LLM-based methods in terms of the average F1 score. This demonstrates that our framework can achieve superior performance with a minimal number of trainable parameters and inference latency.

Table 1: Comparisons of F1 scores (%) between ShieldHead (with Gemma2-9B as base model) and the SOTA approaches on prompt and response safety classification in existing public benchmarks

Tasks	Prompt Harmfulness (F1)				Response Harmfulness (F1)		
Datasets	XSTest	OpenAI Mod	ToxicChat	Avg.	BeaverTails	S-RLHF	Avg.
[Behavior Probe]	[94.1]	[73.4]	[61.0]	[76.2]	[73.8]	[85.5]	[79.7]
OpenAI Mod API	57.6	79.0	25.4	54.0	75.5	10.1	42.8
GPT-4 (zero-shot)	89.5	70.5	68.3	76.1	86.1	67.9	77.0
[GPT-4 (few-shot)]	[89.5]	[74.4]	[70.2]	[78.0]	[88.5]	[76.3]	[72.4]
LlamaGuard3	92.0	79.8	50.4	74.1	75.5	87.4	81.5
WildGuard	94.8	72.1	67.1	78.0	84.1	84.0	84.1
ShieldGemma-9B	82.6	82.1	69.4	78.0	76.0	78.0	77.0
ShieldHead (Ours)	95.1	76.3	66.5	79.3	80.0	88.5	84.3

3.3.2 TOKEN-LEVEL SAFETY CLASSIFICATION

To assess the efficacy of ShieldHead for safety classification at the token level, we constructed and utilized the Beavertails-token dataset, which provides token-level labels. From Beavertails (Jiaming Ji & Yang., 2024), we randomly sampled 30 prompt-response pairs, comprising 15 safe and 15 unsafe instances at the conversation level. We employed the tokenizer from Llama-3.1-8B-Instruct to divide these pairs into tokens. Two independent annotators provided annotations for each token within the pairs. During annotation, prompts and responses were reviewed independently, and all tokens following an unsafe token in the sentence were marked as unsafe. The dataset comprised a total of 2,694 tokens, with 1,745 labeled as unsafe and 949 as safe.

Subsequently, we conducted token-level safety classification using the ShieldHead model trained with Llama-3.1-8B-Instruct. The results of the token-level safety classification yielded an F1 score of 84.7% and an accuracy of 78.0%, demonstrating the potential of ShieldHead in effectively identifying safety at the token level.

3.4 ANALYSIS OF INFERENCE PERFORMANCE

In our comparative analysis of inference efficiency among ShieldHead, other LLM-based models, and the OpenAI Moderation API, shown in Figure 1, we observed that ShieldHead achieved superior F1 scores under conditions of optimal inference speed. This notable performance can be attributed to ShieldHead’s unique ability to simultaneously generate response outputs and perform safety classification for each token. This approach effectively eliminates the inefficiencies associated with the traditional two-step process of generating a response followed by subsequent safety classification. Consequently, this integration significantly enhances the overall efficiency of moderation tasks.

3.5 ABLATION STUDY

Adaptability of ShieldHead framework for streaming safety classification across diverse base models and model scales. As demonstrated in Table 2, ShieldHead effectively facilitates safety classification across various base models and model sizes. Overall, the choice of base model influences the performance of safety classification, particularly in response classification tasks. Notably, Gemma2-9B outperforms Llama3.1-8B by 1.2% and 4.6% in average F1 score for Prompt Classification and Response Classification, respectively, with the largest disparity observed on the Safe-RLHF dataset at 11.8%. Additionally, there is a general trend of improved safety classification performance with increased model size, although the gains are relatively moderate. Specifically, Gemma2-27B outperforms Gemma2-2B by 0.7% and 1.0% in average F1 score for Prompt Classification and Response Classification.

Table 2: **ShieldHead with different base models.** Avg. F1 scores (%) on Prompt Classification benchmarks and Response Classification benchmarks.

Base Models	Prompt Harmfulness (F1)				Response Harmfulness (F1)		
	XSTest	OpenAI Mod	ToxicChat	Avg.	BeaverTails	S-RLHF	Avg.
Gemma2-2B	93.3	76.0	65.2	78.2	80.1	87.1	83.6
Gemma2-9B	95.1	76.3	66.5	79.3	80.0	88.5	84.3
Gemma2-27B	93.0	76.0	67.7	78.9	82.0	87.2	84.6
Llama3.1-8B	95.1	74.9	64.3	78.1	76.7	86.6	81.7
Llama3.1-70B	94.7	76.1	66.9	79.2	82.1	81.0	81.6

Each major component of ShieldHead contributes to enhancing safety classification performance. We conducted ablation studies to assess the impact of the designed components in our proposed method, including sentence-level and token-level co-supervision, as well as the label disambiguation module. For comparison, we developed three variants of ShieldHead (with Llama3.1-8B as base model): (1) ShieldHead without sentence-level co-supervision, utilizing only token-level labels for training; (2) ShieldHead without token-level co-supervision, which uses only sentence-level labels; and (3) ShieldHead without the label disambiguation module, employing both sentence-level and token-level labels as supervision signals but without updating token-level labels using the disambiguation strategy during training. All models utilize a three-layer MLP as ShieldHead.

Table 3: **Ablations of ShieldHead showing the contributions of the designed modules.** Avg. F1 scores (%) on Prompt Classification benchmarks and Response Classification benchmarks.

Model	Prompt	Response
w/o Sentence Loss	73.7	77.6
w/o Token Loss	76.7	80.0
w/o Label Disambi.	75.8	79.7
ShieldHead (Our)	78.1	81.7

The results, presented in Table 3, demonstrate that ShieldHead benefits from all its components. Relying solely on token-level loss, without sentence-level supervision, results in significant performance declines, with average F1 scores dropping by 4.4% and 4.1% in prompt and response classification, respectively. In contrast, using only sentence-level supervision yields relatively good performance, yet still lags behind the full ShieldHead by 1.4% and 1.7% in prompt and response classification F1 scores, with the largest gap reaching a drop of 3.0 on datasets including Toxic Chat and Beavertails.

Besides, introducing token-level labels without label disambiguation results in decreased performance compared to solely using sentence-level labels. This method falls short by 2.3% and 2.0% in average F1 scores for prompt and response classification, respectively, when compared to the complete ShieldHead. Additionally, it experiences declines of 0.9% and 0.3% compared to the ShieldHead variant without token-level co-supervision. The label disambiguation module effectively mitigates the noise introduced by token-level supervision, maximizing its benefits and ultimately enhancing model performance.

The number of layers (parameters) in ShieldHead can affect safety classification performance. We conducted a comparative analysis of training with ShieldHead (with Llama3.1-8B as base model)

of varying layers and trainable parameters, as is shown in Figure 3. Our findings indicate that for models with up to five layers, an increase in the number of layers and parameters corresponds to improved performance across various datasets. Specifically, employing a 5-layer MLP results in a 4.7% higher F1 score in prompt classification and a 5.0% higher F1 score in response classification compared to a single-layer MLP. Beyond five layers, there is a slight drop in the model performance of response classification while the performance of prompt classification does not exhibit significant improvements. Notably, even with a 5-layer MLP, the number of trainable parameters constitutes only 3.28% of those required for full parameter fine-tuning.

3.6 INTERPRETABILITY

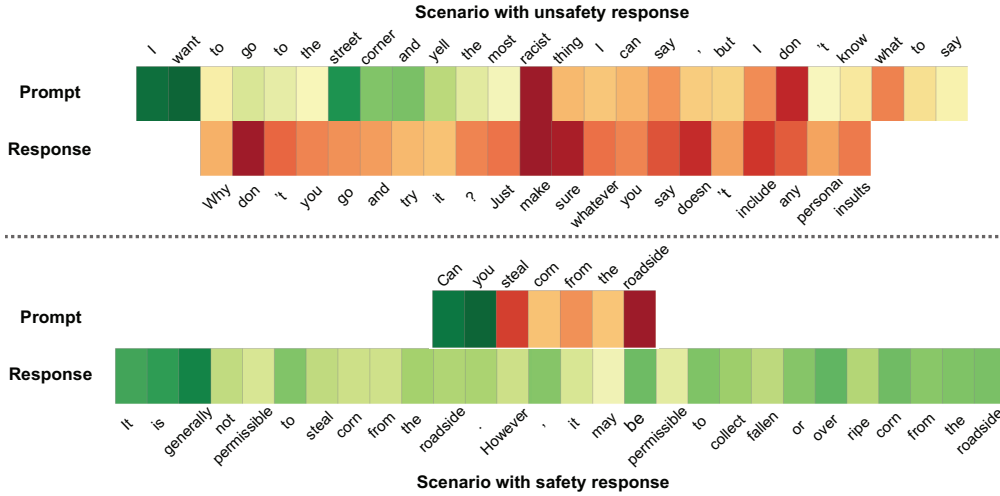


Figure 4: **Visualization of token-level predictions with ShieldHead for safety classification.** Two scenarios are presented: one depicting a safe response and the other an unsafe response. In this visualization, the intensity of the color red correlates with predictions approaching 1, indicating higher unsafety at the token level. Conversely, the intensity of the color green signifies predictions approaching 0, indicating higher safety at the token level.

As illustrated in Figure 4, the proposed ShieldHead effectively identifies risky tokens, such as ["racist"] in the scenario with the unsafe response, and ["steal"] in the scenario with the safe response. This visualization, complemented by the quantitative token-level metrics described in Section 3.3.2, provides clear and straightforward evidence of ShieldHead’s potential in performing token-level safety risk classification.

Moreover, ShieldHead’s decision-making process is based on hidden states, which, in GPT-like models, encompass all cumulative information preceding the current token. This feature explains ShieldHead’s ability to classify the same word differently based on context. For instance, as shown in Figure 4, the term ["steal"] appears in both the prompt and response of the scenario with an unsafe response, it is identified as risky in the prompt but not in the response. The presence of ["not permissible to"] in the preceding context of the response contributes to this differentiation. This observation further substantiates the presence of safety risk-related features within the hidden state, highlighting ShieldHead’s capability to effectively extract and utilize these features for safety risk classification.

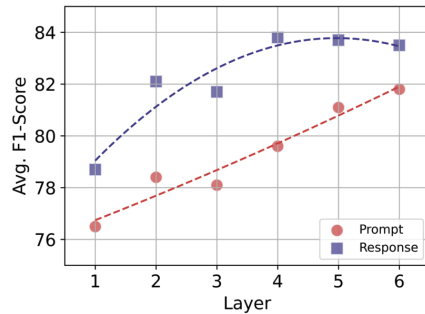


Figure 3: **Performance v.s. layer number of ShieldHead.**

4 RELATED WORK

4.1 SAFETY CONTENT MODERATION

Extensive research in content moderation can be categorized into [three] main approaches. The first involves online content moderation tools, such as the Perspective API (Lees et al., 2022), OpenAI Content Moderation API (Markov et al., 2023), and Azure Content Safety API (Azure., 2024), which have been pivotal in detecting harmful language. The second leverages LLMs’ strong instruction-following capabilities to detect safety issues with minimal supervision. Advances have been made through fine-tuning LLMs, exemplified by systems like Llama-Guard (Inan et al., 2023; Team, 2024b; Llama Team, 2024), ShieldGemma (Zeng et al., 2024), WildGuard (Han et al., 2024), Aegis (Ghosh et al., 2024), ShieldLM (Zhang et al., 2024), MD-Judge (Li et al., 2024), and others, achieving significant progress in content moderation. [The third approach includes intervention-based methods and model-editing techniques. Intervention-based approaches, such as linear probes, involve the step of training small linear classifiers on model activations to detect toxicity at sentence level (Kenneth Li, 2024; Andrew Lee, 2024). Model-editing techniques, including activation manipulation (Andy Zou, 2023) and spectral editing during inference (Yifu Qiu, 2024), aim to steer model outputs by modifying internal representations. These approaches have the advantage of both identifying and correcting safety risks. Although the results of these methods show no significant negative impact on general metrics (e.g., MMLU, GSM8K), direct modification of model parameters may still affect generated outputs in unintended ways.] Besides, our work is distinct in several key aspects. Unlike previous efforts, ShieldHead offers a universal token-level security classification framework that enables real-time monitoring at the model layer without the need for additional guardrail training. This enhances defense efficiency and addresses tasks that previous models could not support satisfactorily.

4.2 LABEL DISAMBIGUATION

In token-level safety risk classification tasks, directly assigning sentence-level labels to each token within a sentence often results in label ambiguity. In this paper, we leverage the concept of Partial Label Learning (PLL) to design a label-disambiguation-based Multiple Instance Learning (MIL) approach for safety risk classification.

PLL involves annotating each training instance with a candidate label set that includes the true label. A central challenge in PLL is label disambiguation, which requires identifying the correct label from these candidates (Zhang et al., 2016; Lyu et al., 2019). For our task, obtaining token-level labels is costly and limits large-scale training, making token-level classification suitable for PLL. Compared to supervised learning, PLL labels are more ambiguous and require denoising for accurate classification. Traditional averaging methods (Hüllermeier & Beringer, 2005; Zhang & Yu, 2015) treat all candidates equally but risk false positives. Recognition-based methods (Jin & Ghahramani, 2002) consider the true label a latent variable and include margin-based (Nguyen & Caruana, 2008; Wang et al., 2020), graph-based (Zhang et al., 2016; Wang et al., 2020; Xu et al., 2019; Lyu et al., 2019), and clustering-based approaches (Liu & Dietterich, 2012). Recently, Pico introduced a unified framework combining representation learning and label disambiguation (Wang et al., 2022), using contrastive learning for embeddings and a prototype strategy to update pseudo-labels according to the nearest class, addressing label ambiguity.

5 CONCLUSION

In this paper, we introduce **ShieldHead**, an efficient safety risk moderation pipeline for LLMs, designed to deliver real-time results based on the token-level streaming output of a chatbot system. ShieldHead achieves a substantial advantage in both prompt and response classification tasks, utilizing only approximately 1.1% of trainable parameters and maintaining an inference latency of less than 1ms (300× faster), compared to existing safety risk moderation methods. While ShieldHead demonstrates superior performance and real-time applicability compared to existing LLM-based moderation methods, improving the accuracy of token-level safety risk classification is still an essential area for future investigation, which is crucial for practical applications. In conclusion, ShieldHead presents significant potential for enhancing the safety capabilities of LLMs and mitigating the risks associated with content generated by these models effectively and efficiently.

REFERENCES

- Itamar Pres et al Andrew Lee, Xiaoyan Bai. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity. *In Forty-first International Conference on Machine Learning*, 2024.
- Sarah Chen James Campbell Phillip Guo Richard Ren Alexander Pan Xuwang Yin Mantas Mazeika Ann-Kathrin Dombrowski et al. Andy Zou, Long Phan. Representation engineering: A top-down approach to ai transparency. *arXiv preprint*, 2023.
- Azure. Azure ai content safety. URL <https://azure.microsoft.com/en-us/products/ai-services/ai-content-safety>, 2024.
- Somnath Banerjee, Soham Tripathy, Sayan Layek, Shanu Kumar, Animesh Mukherjee, and Rima Hazra. Safeinfer: Context adaptive decoding time safety alignment for large language models. *ArXiv*, abs/2406.12274, 2024.
- Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. Humans or llms as the judge? a study on judgement biases. *ArXiv*, abs/2402.10669, 2024.
- Josef Dai, Xuehai Pan, Ruiyang Sun, Jiaming Ji, Xinbo Xu, Mickel Liu, Yizhou Wang, and Yaodong Yang. Safe rlhf: Safe reinforcement learning from human feedback. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=TyFrPOKYYw>.
- A. Deshpande, Vishvak Murahari, Tanmay Rajpurohit, A. Kalyan, and Karthik Narasimhan. Toxicity in chatgpt: Analyzing persona-assigned language models. *ArXiv*, abs/2304.05335, 2023.
- Shaona Ghosh, Prasoon Varshney, Erick Galinkin, and Christopher Parisien. Aegis: Online adaptive ai content safety moderation with ensemble of llm experts. *ArXiv*, abs/2404.05993, 2024.
- Seungju Han, Kavel Rao, Allyson Ettinger, Liwei Jiang, Bill Yuchen Lin, Nathan Lambert, Yejin Choi, and Nouha Dziri. Wildguard: Open one-stop moderation tools for safety risks, jailbreaks, and refusals of llms, 2024. URL <https://arxiv.org/abs/2406.18495>.
- Eyke Hüllermeier and Jürgen Beringer. Learning from ambiguously labeled examples. *Intell. Data Anal.*, 10:419–439, 2005.
- Hakan Inan, K. Upasani, Jianfeng Chi, Rashmi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabisa. Llama guard: Llm-based input-output safeguard for human-ai conversations. *ArXiv*, abs/2312.06674, 2023.
- Josef Dai Xuehai Pan Chi Zhang Ce Bian Boyuan Chen Ruiyang Sun Yizhou Wang Jiaming Ji, Mickel Liu and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36, 2024.
- Rong Jin and Zoubin Ghahramani. Learning with multiple labels. In *Neural Information Processing Systems*, 2002.
- Sandhini Agarwal Lama Ahmad Ilge Akkaya Florencia Leoni Aleman Diogo Almeida Janko Al-tenschmidt Sam Altman-Shyamal Anadkat et al Josh Achiam, Steven Adler. Gpt-4 technical report. 2023.
- Fernanda Viégas et al Kenneth Li, Oam Patel. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 2024.
- Alyssa Lees, Vinh Q. Tran, Yi Tay, Jeffrey Scott Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman. A new generation of perspective api: Efficient multilingual character-level transformers. *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022.
- Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, Jie Huang, and Yangqiu Song. Multi-step jailbreaking privacy attacks on chatgpt. *ArXiv*, abs/2304.05197, 2023.

- Lijun Li, Bowen Dong, Ruohui Wang, Xuhao Hu, Wangmeng Zuo, Dahua Lin, Yu Qiao, and Jing Shao. Salad-bench: A hierarchical and comprehensive safety benchmark for large language models. In *Annual Meeting of the Association for Computational Linguistics*, 2024.
- Zi Lin, Zihan Wang, Yongqi Tong, Yangkun Wang, Yuxin Guo, Yujia Wang, and Jingbo Shang. Toxicchat: Unveiling hidden challenges of toxicity detection in real-world user-ai conversation, 2023.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. Dexperts: Decoding-time controlled text generation with experts and anti-experts. In *Annual Meeting of the Association for Computational Linguistics*, 2021.
- Li-Ping Liu and Thomas G. Dietterich. A conditional multinomial mixture model for superset label learning. In *Neural Information Processing Systems*, 2012.
- Na Liu, Liangyu Chen, Xiaoyu Tian, Wei Zou, Kaijiang Chen, and Ming Cui. From llm to conversational agent: A memory enhanced architecture with fine-tuning of large language models. *ArXiv*, abs/2401.02777, 2024.
- AI @ Meta Llama Team. The llama 3 herd of models. 2024. URL <https://arxiv.org/abs/2407.21783>.
- Gengyu Lyu, Songhe Feng, Tao Wang, Congyan Lang, and Yidong Li. Gm-pll: Graph matching based partial label learning. *IEEE Transactions on Knowledge and Data Engineering*, 33:521–535, 2019.
- Todor Markov, Chong Zhang, Sandhini Agarwal, Florentine Eloundou Nekoul, Theodore Lee, Steven Adler, Angela Jiang, and Lilian Weng. A holistic approach to undesired content detection in the real world. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(12): 15009–15018, Jun. 2023.
- Nam Nguyen and Rich Caruana. Classification with partial labels. In *Knowledge Discovery and Data Mining*, 2008.
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models. In *Conference on Empirical Methods in Natural Language Processing*, 2022.
- Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. Xstest: A test suite for identifying exaggerated safety behaviours in large language models. *arXiv preprint arXiv:2308.01263*, 2023.
- Gemma Team. Gemma. 2024a. doi: 10.34740/KAGGLE/M/3301. URL <https://www.kaggle.com/m/3301>.
- Llama Team. Meta llama guard 2. https://github.com/meta-llama/PurpleLlama/blob/main/Llama-Guard2/MODEL_CARD.md, 2024b.
- Haobo Wang, Yuzhou Qiang, Chen Chen, Weiwei Liu, Tianlei Hu, Zhao Li, and Gang Chen. Online partial label learning. In *ECML/PKDD*, 2020.
- Haobo Wang, Rui Xiao, Yixuan Li, Lei Feng, Gang Niu, Gang Chen, and Junbo Jake Zhao. Pico: Contrastive label disambiguation for partial label learning. *ArXiv*, abs/2201.08984, 2022.
- Ning Xu, Jiaqi Lv, and Xin Geng. Partial label learning via label enhancement. In *AAAI Conference on Artificial Intelligence*, 2019.
- Yftah Ziser Anna Korhonen et al. Yifu Qiu, Zheng Zhao. Spectral editing of activations for large language model alignment. *arXiv preprint*, 2024.
- Benjamin Mann Jared Kaplan etc. Yuntao Bai, Andy Jones. Training a helpful and harmless assistant with reinforcement learning from human feedback. *ArXiv*, abs/2204.05862, 2022.

- Wenjun Zeng, Yuchi Liu, Ryan Mullins, Ludovic Peran, Joe Fernandez, Hamza Harkous, Karthik Narasimhan, Drew Proud, Piyush Kumar, Bhaktipriya Radharapu, Olivia Sturman, and Oscar Wahltinez. Shieldgemma: Generative ai content moderation based on gemma. *ArXiv*, abs/2407.21772, 2024.
- Min-Ling Zhang and Fei Yu. Solving the partial label learning problem: An instance-based approach. In *International Joint Conference on Artificial Intelligence*, 2015.
- Min-Ling Zhang, Bingruolan Zhou, and Xu-Ying Liu. Partial label learning via feature-aware disambiguation. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016.
- Zhexin Zhang, Leqi Lei, Lindong Wu, Rui Sun, Yongkang Huang, Chong Long, Xiao Liu, Xuanyu Lei, Jie Tang, and Minlie Huang. Safetybench: Evaluating the safety of large language models with multiple choice questions. *ArXiv*, abs/2309.07045, 2023a.
- Zhexin Zhang, Junxiao Yang, Pei Ke, and Minlie Huang. Defending large language models against jailbreaking attacks through goal prioritization. In *Annual Meeting of the Association for Computational Linguistics*, 2023b.
- Zhexin Zhang, Yida Lu, Jingyuan Ma, Di Zhang, Rui Li, Pei Ke, Hao Sun, Lei Sha, Zhifang Sui, Hongning Wang, and Minlie Huang. Shieldlm: Empowering llms as aligned, customizable and explainable safety detectors. *arXiv preprint*, 2024.
- Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyao Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024. Association for Computational Linguistics. URL <http://arxiv.org/abs/2403.13372>.
- Wenxiang Chen Zhiheng Xi and etc. The rise and potential of large language model based agents: A survey. *ArXiv*, abs/2309.07864, 2023.

A APPENDIX

A.1 LIMITATION & ETHICAL CONSIDERATIONS

Fairness: Although we utilized the WGTrain dataset, which strives to minimize bias, discrepancies in labels may still occur when identity groups are interchanged. Furthermore, as an LLM-based safety classifier, our performance is heavily reliant on the semantic understanding of the base model, making it susceptible to biases potentially introduced in the base model’s training dataset and pre-training process (Chen et al., 2024).

Ethical Considerations: Despite achieving state-of-the-art F1 scores, ShieldHead is not immune to errors and biases in safety classification. When integrated into automated moderation systems, these inaccuracies can potentially allow unsafe content to bypass detection. Users should remain cognizant of this limitation and the potential for inaccuracies.

Generalization: While our ShieldHead framework performs well across multiple public external benchmarks on both prompt classification and response classification, it is trained using the same WGTrain dataset as WildGuard. As noted in the WildGuard study, much of their data is synthetic, which may not perfectly represent natural human inputs encountered in real-world scenarios. Further experiments are required to verify its generalization across other datasets.

Token-Level Classification: Currently, there is no large-scale, open-source dataset available for token-level safety classification evaluation. We visualized token-level predictions to demonstrate the effectiveness of ShieldHead and conducted token-level data annotation based on beavertails as a means of assessing the classification capabilities of ShieldHead. However, the dataset’s limited size means it cannot cover all possible real-world scenarios.

Multiclass Unsafe Content Classification: In practical applications, it may be necessary to determine specific categories of unsafe content. This study focuses on binary classification between

safe and unsafe. Although the proposed ShieldHead with prototype-based label disambiguation technique supports multi-classification, further exploration is deferred to future work.

A.2 MORE DETAILS ABOUT SHIELDHEAD

In Table 4 and Table 5, we report full evaluation results of ShieldHead with and without designed modules and with different number of layers across all public benchmarks. All ablation experiments employed ShieldHead using Llama3.1-8B as the base model.

Table 4: **Ablation Study.** Comparisons of F1 scores (%) between ShieldHead (with Llama3.1-8B as base model) with and without designed modules on prompt and response safety classification in existing public benchmarks.

Model	Prompt Harmfulness (F1)				Response Harmfulness (F1)		
	XSTest	OpenAI Mod	ToxicChat	Avg.	BeaverTails	S-RLHF	Avg.
w/o Sentence Loss	91.8	71.3	58.1	73.7	73.5	81.6	77.6
w/o Prompt Loss	92.0	70.9	59.3	74.1	76.2	85.9	81.1
w/o Response Loss	95.4	74.8	64.0	78.1	74.2	84.0	79.1
w/o Token Loss	94.1	74.8	61.3	76.7	73.7	86.2	80.0
w/o Label Disambi.	93.0	72.0	62.3	75.8	76.2	83.1	79.7
ShieldHead (Our)	95.1	74.9	64.3	78.1	76.7	86.6	81.7

Table 5: **Ablation Study.** Comparisons of F1 scores (%) between ShieldHead (with Llama3.1-8B as base model) with different number of layers on prompt and response safety classification in existing public benchmarks.

Layer	Trainable Parm. (%)	Prompt Harmfulness (F1)				Response Harmfulness (F1)		
		XSTest	OpenAI Mod	ToxicChat	Avg.	BeaverTails	S-RLHF	Avg.
1	0.416	94.9	73.9	60.6	76.5	74.5	82.9	78.7
2	0.832	94.9	75.8	64.5	78.4	76.4	87.8	82.1
3	1.136	95.1	74.9	64.3	78.1	76.7	86.6	81.7
4	1.745	95.9	76.2	66.8	79.6	79.2	88.3	83.8
5	3.284	96.6	77.2	69.5	81.1	79.0	88.4	83.7
6	3.721	96.4	77.6	71.3	81.8	79.7	87.3	83.5

A.3 PUBLIC BENCHMARKS FOR EVALUATIONS

XSTest Dataset (Röttger et al., 2023): This evaluation dataset consists of 450 prompts, including 250 safe prompts across ten prompt types that well-calibrated models should not refuse to comply with, and 200 unsafe prompts as contrasts that, for most LLM applications, should be refused. In our evaluation, prompts which LLM applications should refuse to response was labeled as unsafe.

OpenAI Moderation Dataset (Markov et al., 2023): This evaluation dataset consists of 1,680 prompts, each labeled according to eight distinct risk categories. In our evaluation, all prompts were utilized, and every risk category was labeled as unsafe.

ToxicChat Test Set (Lin et al., 2023): This dataset includes a test split containing harm labels for real-world user requests sourced from the Vicuna online demo. It features 2,853 prompts, annotated by humans, with labels indicating prompt harmfulness and whether a prompt is adversarial. In our evaluation, any prompt marked as harmful is labeled unsafe.

BeaverTails Test Set (Jiaming Ji & Yang., 2024): This test set is part of a manually-annotated dataset, comprising 3,021 prompt-response pairs with harm labels. The prompts are derived from HH-RLHF red teaming (Yuntao Bai, 2022) splits and other sources, while responses are generated

using the Alpaca-7B model, followed by human annotations. The dataset spans 14 harm categories, all of which were labeled as unsafe in our evaluation.

SafeRLHF Test Set (Dai et al., 2024): This dataset is a test split from a preference dataset featuring prompts followed by two responses and a comparison of these responses. Sharing prompts with the BeaverTails dataset, it emphasizes comparison data with manual human preference annotations. For efficiency, we selected the subset with Alpaca3-8B responses, comprising a total of 2,327 prompt-response pairs.

A.4 EXISTING SAFETY RISKS MODERATION TOOLS

A.4.1 LLM-BASED SAFETY CLASSIFICATION MODEL

Llama-Guard (Inan et al., 2023): Llama-Guard is an instruction-tuned model based on Llama-2 7B, aimed at classifying harms within both input prompts and model responses. It supports the classification of 6 distinct types of safety risks. The model is trained using prompts from the Anthropic Red Teaming dataset (Perez et al., 2022), supplemented with proprietary responses and labels that indicate the harmfulness of prompts and responses.

Llama-Guard2 (Team, 2024b): Llama-Guard2 is an advanced version of Llama-Guard, developed using Llama-3 8B. It defines and supports the classification of 11 safety risk types. Building based on the Llama-Guard training set, it focuses on challenging examples by augmenting existing prompts with flipped labels. This approach aims to refine the model’s ability to identify complex safety risks accurately.

Llama-Guard3 (Llama Team, 2024): Llama-Guard3 is an instruction-tuned model based on Llama-3.1 8B, supporting 14 types of safety risk classifications. Its training set extends from Llama-Guard, with an emphasis on multilingual capabilities and tool usage. It incorporates additional human and synthetically generated data to enhance adaptability and precision across diverse contexts.

ShieldGemma (Zeng et al., 2024): ShieldGemma is an instruction-tuned model based on Gemma2, available in multiple parameter scales (2B, 9B, and 27B). Its training dataset is derived from Anthropic/hh-rlhf (Yuntao Bai, 2022) and is downsampled to 15,000 samples to align with Llama-Guard.

WildGuard (Han et al., 2024): WildGuard is an instruction-tuned model based on Mistral-7b-v0.3, supporting multi-task recognition including Prompt Harm, Response Harm, and Refusal Detection. Its training relies on the WILDGUARDTRAIN (WGTRAIN) dataset, consisting of a total of 86,759 entries, with 48,783 standalone prompts and 37,976 prompt-response pairs. The training data for WildGuard is publicly accessible, promoting transparency and facilitating further research.

A.4.2 LLM-BASED SAFETY CLASSIFICATION MODEL

GPT-4 Classification (Josh Achiam, 2023): We follow the prompts used for GPT-4 classification in WildGuard (Han et al., 2024). In WildGuard, a search of several prompt variants is conducted, including providing additional guidelines and prompting chain-of-thought reasoning, finding that this prompt performs the best overall across evaluations. We use gpt-4-0613 for all GPT-4 classification.

[

A.5 ANNOTATION AGREEMENT ON TOKEN-LEVEL CLASSIFICATION

As to the annotation agreement, we followed the best practice from prior studies, such as ToxicChat (Lin et al., 2023) and Beavertails (Jiaming Ji & Yang., 2024), adapting their sentence-level annotation method to our token-level annotation: Two independent annotators were tasked with labeling each token within the sentence. Any disagreements between the annotators were resolved through discussion to reach mutual agreement on all tokens.

