

ToxiTrace: Gradient-Aligned Training for Explainable Chinese Toxicity Detection

Anonymous ACL submission

Abstract

Existing Chinese toxic content detection methods mainly target sentence-level classification but often fail to provide readable and contiguous toxic evidence spans. We propose **ToxiTrace**, an explainability-oriented method for BERT-style encoders with three components: (1) **CuSA**, which refines encoder-derived saliency cues into fine-grained toxic spans with lightweight LLM guidance; (2) **GCloss**, a gradient-constrained objective that concentrates token-level saliency on toxic evidence while suppressing irrelevant activations; and (3) **ARCL**, which constructs sample-specific contrastive reasoning pairs to sharpen the semantic boundary between toxic and non-toxic content. Experiments show that ToxiTrace improves classification accuracy and toxic span extraction while preserving efficient encoder-based inference and producing more coherent, human-readable explanations.

Disclaimer: *This paper describes violent and discriminatory content that maybe disturbing to some readers.*

1 Introduction

In the era of pervasive digital social media, toxic user-generated content (UGC)—such as cyberbullying and hate speech—has become increasingly prevalent, posing tangible risks to online communities and society at large. As a result, toxic content detection has been extensively studied (Arora et al., 2023; Kirk et al., 2022; Azumah et al., 2023), with Transformer-based pre-trained language models (PLMs) (Vaswani et al., 2017; Devlin et al., 2019; Liu et al., 2019) and, more recently, large language models (LLMs) further advancing classification performance.

Despite these advances, most existing work focuses on sentence-level toxicity classification, while providing little insight into which specific

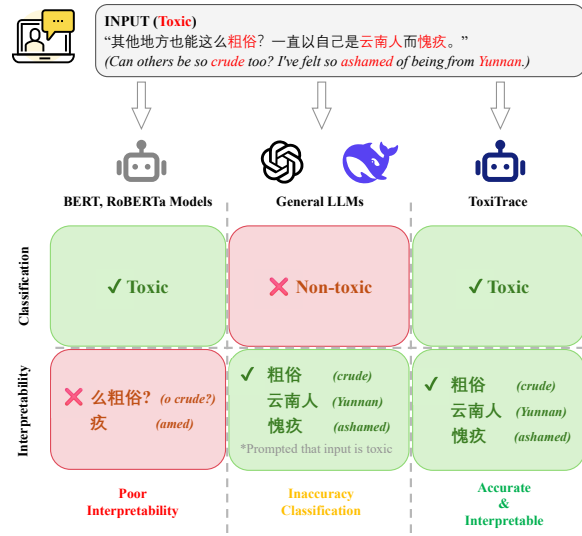


Figure 1: Existing encoder-based detectors struggle to reliably extract fine-grained toxic expressions within a sentence; LLMs can better extract spans when toxicity is given but are limited in direct classification and efficiency; our method preserves classification performance while enabling contiguous toxic span extraction.

parts of a sentence constitute the toxic content. However, identifying fine-grained toxic evidence is crucial for explainability, moderation transparency, and downstream interventions. This challenge is particularly pronounced for Chinese.

Unlike English, where words serve as the basic semantic units, Chinese toxic expressions are typically realized as multi-character phrases, while individual characters are often semantically ambiguous. However, mainstream Chinese PLMs adopt character-level tokenization (Cui et al., 2020; Sun et al., 2021). Consequently, attribution signals such as gradients or attention weights are fragmented across individual characters rather than coherent semantic spans, producing rationales that are difficult for humans to interpret (Ding and Koehn, 2021).

As a result, existing Chinese toxic con-

059 tent detection approaches—despite achieving
060 strong sentence-level performance through fine-
061 tuning (Deng et al., 2022), knowledge distilla-
062 tion (Deng et al., 2023), or glyph-aware model-
063 ing (Wullach et al., 2022; Xue et al., 2025)—
064 remain limited in their ability to accurately ex-
065 tract the true toxic expressions within a sentence
066 (Figure 1, left). In contrast, LLMs often exhibit
067 stronger capabilities in explanation and span ex-
068 traction (Creswell and Shanahan, 2022; Chuang
069 et al., 2025), but they typically underperform on
070 direct toxicity classification and incur substan-
071 tially higher inference costs (Schmidhuber and Kr-
072 uschwitz, 2024; Zhang et al., 2025; Sun et al.,
073 2023) (Figure 1, middle).

074 These limitations call for a method that pre-
075 serves the classification strength and efficiency
076 of encoder-based models while enabling reliable
077 extraction of contiguous, human-readable toxic
078 spans (Figure 1, right). To this end, we propose
079 **ToxiTrace**, a span-extraction-oriented framework
080 built on BERT-style encoders for Chinese toxic
081 content detection, designed to produce coherent
082 and interpretable toxic rationales without requir-
083 ing fine-grained span supervision. Our main con-
084 tributions are as follows:

- 085 • We propose a cue-guided span annotation
086 strategy with gradient-aware training, which
087 leverages attribution signals from encoder
088 models to induce consistent saliency on toxic
089 tokens without requiring explicit span anno-
090 tations.
- 091 • We introduce a bidirectional cliff-based span
092 extraction algorithm that identifies contigu-
093 ous toxic spans based on saliency transitions,
094 effectively alleviating the span fragmentation
095 issue inherent in prior top- k selection meth-
096 ods.
- 097 • We develop an adversarial reasoning con-
098 trastive learning objective with adaptive In-
099 foNCE, which aligns sample-specific toxic
100 and non-toxic reasoning representations,
101 sharpening the semantic boundary and fur-
102 ther enhancing span-level interpretability.
- 103 • Extensive experiments on multiple Chinese
104 toxic content benchmarks demonstrate that
105 **ToxiTrace** consistently outperforms strong
106 baselines.

2 Related Work 107

2.1 Toxic Content Detection 108

109 Toxic content detection has evolved from early
110 Bag-of-Words and conventional classifiers (Kwok
111 and Wang, 2013; Waseem and Hovy, 2016; David-
112 son et al., 2017) to neural and Transformer-based
113 models that achieve strong sentence-level accu-
114 racy (Badjatiya et al., 2017; Zimmerman et al.,
115 2018; Caselli et al., 2021; Sarkar et al., 2021).

116 Chinese toxic content detection follows a sim-
117 ilar trajectory, supported by datasets such as
118 COLD (Deng et al., 2022), ToxiCN (Lu et al.,
119 2023), and CNTP (Yang et al., 2025), as well as
120 encoder-centric methods including character-level
121 modeling (Wullach et al., 2022), domain feature
122 fusion (Zhang et al., 2024), LLM-assisted rewrit-
123 ing (Chao et al., 2024), and distillation for robust-
124 ness (Deng et al., 2023). In contrast to prior work
125 that primarily optimizes sentence-level detection,
126 we target fine-grained toxic span extraction by
127 training an encoder to produce evidence-aligned
128 saliency and readable spans.

129 Although CNTP (Yang et al., 2025) provides
130 limited span annotations, most Chinese detection
131 pipelines still lack reliable supervision and meth-
132 ods for extracting *contiguous* toxic spans within
133 sentences. This gap often leaves models accurate
134 yet poorly grounded in human-readable evidence.
135 We address this by using cue-guided span signals
136 to support training under weak supervision and by
137 enforcing higher gradient responses on toxic evi-
138 dence.

2.2 Attribution Method 139

140 Attribution methods, initially developed in com-
141 puter vision, have been widely adopted for
142 NLP interpretability. Representative lines in-
143 clude perturbation-based explanations such as
144 LIME (Ribeiro et al., 2016), gradient-based ex-
145 planations (Ross et al., 2017), and task-calibrated
146 attention-saliency alignment that improves faith-
147 fulness (Chrysostomou and Aletras, 2021a,b). A
148 practical issue is that many pipelines extract ex-
149 planations by selecting top- k salient tokens, where
150 k is either fixed (Jesus et al., 2021; Bastings et al.,
151 2022) or set by heuristics such as a fixed frac-
152 tion of sentence length (Krishna et al., 2025); dy-
153 namic alternatives such as peak-based top- k have
154 been proposed to improve consistency (Kamp
155 et al., 2023). Building on gradient-based attribu-
156 tion, we explicitly *train* the encoder to yield more

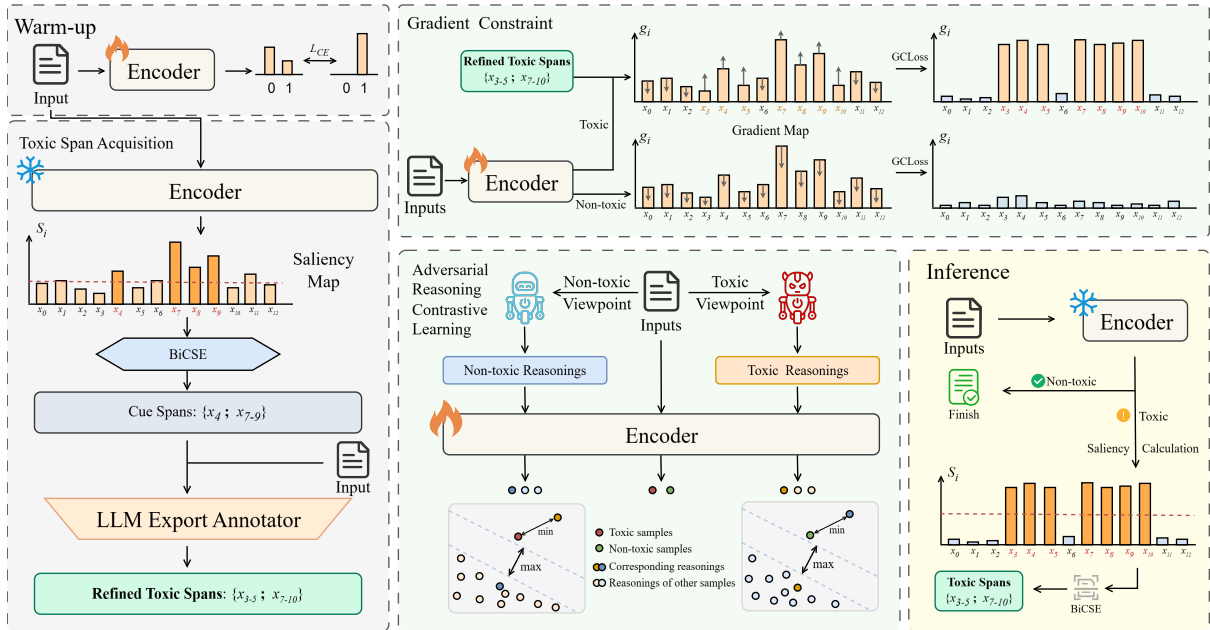


Figure 2: Framework of the proposed ToxiTrace method. During training, we warm up an encoder classifier, acquire weak span annotations with CuSA using BiCSE-extracted saliency cues, and jointly optimize GCLoss and ARCL to concentrate saliency on toxic evidence. During inference, the model predicts toxicity and, for toxic inputs, extracts contiguous spans via BiCSE from the saliency map.

evidence-aligned gradients, rather than only post-hoc selecting salient tokens.

Despite progress, top- k selection often produces fragmented highlights and fails to recover contiguous, human-readable spans—an issue that is especially pronounced under character-level tokenization. We therefore propose a bidirectional scanning algorithm that identifies consecutive locally high-saliency spans, enabling stable contiguous toxic span extraction beyond discrete token selection.

3 Methodology

As shown in Figure 2, our ToxiTrace framework contains following four steps: (1) We first warm up the encoder with standard classification training to obtain robust sentence-level discrimination. (2) After warming-up, we compute saliency maps and apply a **B**idirectional **C**liff-based **S**pan **E**xtraction algorithm (BiCSE) to obtain initial high-saliency spans, which are used as cues to prompt an LLM to refine boundaries and recover coherent toxic spans, yielding *refined toxic spans* as weak span annotations. (3) We then introduce **G**radient **C**onstraint **L**oss (GCLoss) to explicitly increase gradient responses on toxic evidence while suppressing spurious activations on non-toxic tokens, shaping a more concentrated and extractable saliency. (4) In paral-

lel, we adopt **A**dversarial **R**easoning **C**ontrastive **L**earning (ARCL) with adaptive InfoNCE to align each input with sample-specific reasoning of opposing stances, sharpening the toxic/non-toxic semantic boundary.

At inference time, the model first predicts toxicity; if toxic, the saliency map will be calculated, and toxic spans will be extracted using BiCSE.

3.1 Cue-guided Span Annotation

Exist toxic content datasets only have coarse-grained labels (toxic/non-toxic) and cannot accurately locate toxic spans. CuSA constructs span-level signals by using the model’s attribution map as cues and letting an LLM refine span boundaries. It consists of two steps: (1) warm-up fine-tuning to obtain a reliable sentence-level classifier; and (2) cue-guided span refinement, where we feed the toxic text together with initially extracted salient spans as cues to an LLM for span annotation.

Warm-up training. We train the encoder with binary cross-entropy. Given an input text sequence $\mathcal{X} = \{x_1, \dots, x_n\}$, the embedding layer maps tokens to $\mathcal{E} = [e_1, \dots, e_n]$, the model ϕ generate its contextual representation $\mathcal{H} = \phi(\mathcal{E}) = [\mathbf{h}^{cls}, \mathbf{h}^1, \mathbf{h}^2, \dots, \mathbf{h}^n]$. This representation is then fed into a classification head ψ to yield the predicted probability $P(y|\mathcal{X}) = \psi(\mathbf{h}^{cls})$.

Saliency cues for span annotation. After warm-up, following exists attribution methods (Chrysostomou and Aletras, 2021a; Sikdar et al., 2021), we compute a token-level *saliency* score for each token and form a saliency map, which serves as cues for span extraction. The saliency s_i of each token x_i can be calculated via Eq. (1):

$$s_i = \left\| e_i \odot \frac{\partial \log P(y|\mathcal{X})}{\partial e_i} \right\|_2. \quad (1)$$

Exist attribution methods select the top- k most salient tokens (Kamp et al., 2023); however, these selected tokens tend to be scattered. To address this limitation, we propose BiCSE (detailed in Appendix A), a bidirectional cliff-based scanning algorithm that tracks saliency transitions to identify the optimal start and end boundaries of contiguous spans. Moreover, we capture longer and more continuous toxic spans by taking the union of results from two sequential scans (left-to-right and right-to-left).

To help the LLM to find more potential toxic spans, both the raw text and the initially extracted spans (cues) \mathcal{T}_{cue} are fed into the LLM. In our experiments, we use Gemini 2.5 Pro (Team, 2025) as an expert annotator to integrate and refine these cues. The refined toxic spans are formulated as: $\mathcal{T}_{refined} = \text{LLM}(\mathbf{x}, \mathcal{T}_{cue})$. By leveraging the LLM’s superior interpretability, we can achieve more accurate annotation of toxic spans.

3.2 Gradient Constraint Loss

The fine-tuned model attains satisfactory classification performance overall; yet, its token-level toxicity discrimination remains imprecise (Appendix B shows the saliency maps before an after applying our ToxiTrace method). CuSA provides refined toxic spans as annotations, which allow us to explicitly shape the model’s token-level attribution for span extraction. Concretely, GCLOSS consists of two complementary components: (i) a **Pairwise Gradient Ranking (PGR)** term that enforces a margin between toxic and non-toxic tokens, and (ii) a **Push-Pull Threshold (PPT)** term that regularizes their gradient ranges within each sample. During training, both terms operate on the gradient norm of the log-predicted probability with respect to the *input embeddings* e_i :

$$g_i = \left\| \frac{\partial \log P(y|\mathcal{X})}{\partial e_i} \right\|_2 \quad (2)$$

as the training signal to increase responses on toxic spans and suppress spurious activations on non-toxic tokens, thereby shaping more concentrated and extractable attribution.

PGR Loss. The objective of this loss is to penalize cases where, within a single sentence, the gradient of a toxic token exceeds that of a non-toxic token by a margin smaller than the predefined threshold m . The specific formula is given as follows:

$$\mathcal{L}_{PGR} = \frac{1}{|\mathcal{P}| \cdot |\mathcal{N}|} \sum_{i \in \mathcal{P}} \sum_{j \in \mathcal{N}} \max(0, g_j - g_i + m), \quad (3)$$

where \mathcal{P} and \mathcal{N} denote the sets of toxic and non-toxic tokens, respectively; g_i, g_j is calculated according to Eq. (2); and m is set to 1 in this study.

PPT Loss. The PGR Loss introduced above captures the relative gradient relationship between toxic and non-toxic tokens, yet it cannot constrain the gradient value ranges of either token type. Given the substantial discrepancies in gradient scores across different sentences, employing a fixed range to regulate the gradients of toxic and non-toxic tokens is inherently flawed. To address this dual limitation, we propose the intra-sentence PPT Loss that leverages gradient information of tokens within each single sample to separately guide the gradient learning of toxic and non-toxic tokens. Specifically, we first calculate the 15th percentile of the gradient values of all tokens in a single sample as the threshold τ , which serves to constrain the gradient values of non-toxic tokens to stay below this threshold. The detailed formula is given as follows:

$$\mathcal{L}_{neg} = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}} [\max(0, g_j - \tau)]. \quad (4)$$

For toxic tokens, we expect their gradient values to fall within a relatively high range. Thus, we adopt the maximum gradient g_{\max} as the reference and set $\alpha \cdot g_{\max}$ as the lower bound, where α denotes a positive target coefficient. Our goal is to push the gradient values of toxic tokens above this threshold. However, given the possibility of an excessively large g_{\max} , we further introduce a gradient cap τ_{cap} to prevent gradient explosion caused by overly high gradient values of toxic tokens. The formula for toxic tokens is given as follows:

$$\mathcal{L}_{pos} = \frac{1}{|\mathcal{P}|} \sum_{i \in \mathcal{P}} [\max(0, t_n - g_i)], \quad (5)$$

where the target value $t_n = \min(\alpha \cdot g_{max}, \tau_{cap})$. The overall form of the PPT Loss is:

$$\mathcal{L}_{PPT} = \frac{1}{2}(\mathcal{L}_{pos} + \mathcal{L}_{neg}) \quad (6)$$

Remark. GCLOSS constrains *gradients* during training (Eq. (2)), because gradients directly capture the model’s sensitivity to token-level evidence. For span extraction at inference time, we compute a *saliency* score with the same embedding-level gradients.

This formulation reflects both sensitivity and contribution, and BiCSE is applied to the saliency sequence $\{s_i\}_{i=1}^n$ to produce contiguous spans.

3.3 Adversarial Reasoning Contrastive Learning

While the aforementioned loss functions address token-level gradient constraints, they are confined to individual sentences and cannot capture the semantic boundary differentiating toxic from non-toxic sentences. Inspired by the work of (Rusak et al., 2025), we adopt an adaptive InfoNCE loss to implement semantic contrastive learning. Existing data augmentation methods (Zhou et al., 2021; Hu et al., 2023; Fang et al., 2024) primarily rely on substituting, modifying or add noise to certain words in a sentence. Such operations lack a thorough understanding of the sentence’s inherent semantics. To address this limitation, we propose leveraging the **LLM debate mechanism** (Moniri et al., 2025) to generate **adversarial reasoning** as augmented samples. This approach enables the model to better explore the intrinsic semantic information of sentences and learn to distinguish the differences between toxic and non-toxic texts and sharpening the semantic boundary.

Specifically, we first use two opposing reasoning prompts: "Assuming the text is {toxic, normal}, generate explanations to support this judgment"¹, and feed them to the LLM (Gemini 2.5 Flash) to obtain semantically augmented positive and negative samples. This way, the model can produce targeted reasoning content, instead of generating generic descriptions such as "This sentence is just an exaggerated joke and does not intend to attack any group".

¹Detailed prompt templates are provided in Appendix C.

Using the two opposing reasoning prompts obtained above, we adopt the following contrastive loss to facilitate sentence-level semantic learning:

$$\mathcal{L}_{\{tox,nor\}} = -\log \frac{\exp(\mathbf{h}_t \cdot \mathbf{h}_{\{p,n\}}/\tau)}{\exp(\mathbf{h}_t \cdot \mathbf{h}_{\{p,n\}}/\tau) + \sum_{\mathbf{h}_{\{n,p\}}^k \in \mathcal{B}} \exp(\mathbf{h}_t \cdot \mathbf{h}_{\{n,p\}}^k/\tau)}, \quad (7)$$

where \mathbf{h}_t , \mathbf{h}_p and \mathbf{h}_n denote the semantic embeddings of the target text, its positive reasoning explanation and its negative reasoning explanation, respectively; \mathbf{h}_p^k and \mathbf{h}_n^k denote the final-layer [CLS] token embeddings of the positive and negative reasoning explanations corresponding to other toxic and non-toxic sentences within the training batch \mathcal{B} . τ denotes the temperature parameter, which is set to 0.05 in our experiments.

The overall semantic contrastive loss is defined as the average of these two loss components:

$$\mathcal{L}_{con} = \frac{1}{2}(\mathcal{L}_{tox} + \mathcal{L}_{nor}) \quad (8)$$

3.4 Joint Training Objective

In the joint training phase, we simultaneously optimize three loss components: the gradient-based binary classification loss, the gradient constraint losses, and the contrastive learning loss. The overall training objective is formulated as:

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda_{grad}(\mathcal{L}_{PGR} + \mathcal{L}_{PPT}) + \lambda_{sem}\mathcal{L}_{con}, \quad (9)$$

where λ_{grad} and λ_{sem} are hyperparameters that balance the contributions of the gradient constraint losses and the semantic contrastive loss, respectively.

The overall training pipeline is designed as follows: 1. First, perform a warm-up training using the cross-entropy loss to equip the model with basic toxicity classification capability. 2. Subsequently, the GCLOSS is introduced to enforce the model to generate higher gradient values for toxic tokens and lower gradient values for non-toxic tokens. 3. Meanwhile, ARCL is adopted to sharpen the semantic boundary between toxic and non-toxic texts.

4 Experiments

4.1 Experimental Setup

Dataset. For the binary toxicity classification task, we evaluate our model on two Chinese

Models	COLD					ToxiCN				
	<i>Acc</i>	<i>R</i>	<i>P</i>	<i>F</i> ₁	Macro- <i>F</i> ₁	<i>Acc</i>	<i>R</i>	<i>P</i>	<i>F</i> ₁	Macro- <i>F</i> ₁
Qwen3-8B	74.03	59.90	70.15	64.62	72.33	71.44	82.94	69.13	75.41	71.51
LLaMA3.1-8B	73.91	51.46	74.58	60.90	72.00	65.18	44.29	81.61	57.42	68.25
DeepSeek-V3.2	74.23	59.23	70.91	64.55	72.51	60.66	36.19	77.35	49.30	64.14
GPT-4o	74.49	67.79	66.98	67.38	73.22	73.20	65.66	66.98	67.38	73.22
RoBERTa	82.75	86.43	74.24	79.87	82.76	82.70	84.38	83.40	83.89	82.81
Qwen3-8B (SFT)	82.25	82.68	75.02	78.66	81.87	82.08	83.52	82.74	83.12	82.02
LLaMA2-7B (SFT)	82.43	86.28	73.78	79.54	82.46	81.25	84.62	80.81	82.67	81.18
LLaMA3.1-8B (SFT)	83.00	<u>87.13</u>	74.37	80.19	82.19	83.02	86.08	82.52	84.26	82.96
LLaMA3.1-8B + ours	82.15	88.42	72.52	79.68	81.89	82.33	84.93	82.22	83.55	82.23
Qwen3-8B + ours	82.96	82.53	<u>76.34</u>	79.32	82.41	82.91	83.44	84.10	83.77	82.86
RoBERTa + ours	83.84	86.19	76.14	80.85	83.68	<u>83.62</u>	<u>87.05</u>	<u>82.82</u>	84.88	83.56
MacBERT + ours	<u>83.22</u>	86.19	<u>75.10</u>	<u>80.27</u>	<u>83.13</u>	83.87	87.99	82.61	85.21	83.83

Table 1: The average classification results of different models on COLD and Toxicn datasets(%) across 3 independent runs.

datasets: COLD (Deng et al., 2022) (32,480 instances in total, with 5,323 test instances) and ToxiCN (Lu et al., 2023) (12,011 instances in total, with 2,411 test instances).

For toxic span extraction, we use the span annotations provided in CNTP (Yang et al., 2025) as gold labels, which include 2,533 samples annotated with toxic spans, to assess the model’s ability to extract fine-grained toxic expressions within toxic sentences.

Models. To verify the effectiveness of our proposed method in both binary classification and span-level extraction, we applied our training strategy across different BERT-based models: **Chinese-RoBERTa-wwm-ext** (RoBERTa) and **Chinese-MacBERT-base** (MacBERT) (Cui et al., 2020). Concurrently, a set of LLMs were chosen to conduct comparative analysis: **LLaMA2-7B**, **LLaMA3.1-8B**, **Qwen3-8B**, **DeepSeek-V3** and **GPT-4o**.

Implement. Corresponding models fine-tuned solely on binary classification labels were set as baselines. For open-source LLMs, both their direct inference capabilities and their performance after fine-tuning on the respective datasets (denoted as SFT) were evaluated. Fine-tuning of LLMs was conducted using the open-source toolkit LLaMA-Factory² to implement LoRA (Hu et al., 2022). Closed-source LLMs were evaluated in a zero-shot setting; the prompt templates used were detailed in Appendix C. Models denoted with "(ours)" were trained using the method

²<https://github.com/hiyouga/LLaMA-Factory>

proposed in this paper. DeepSeek and GPT-4o were accessed via official APIs, while all other experiments were conducted using four NVIDIA A800 80GB GPUs.

Evaluation metrics. To evaluate toxic content detection performance, we used five widely adopted metrics: Accuracy (*Acc*), Recall (*R*), Precision (*P*), *F*₁ and Macro-*F*₁ Score.

4.2 Overall Classification Performance

The classification performance of ToxiTrace and competing methods is summarized in Table 1. Overall, ToxiTrace achieves the best results across all five metrics on both COLD and ToxiCN. In the zero-shot setting, both open- and closed-source LLMs perform poorly on Chinese toxicity classification (and LLaMA2-7B fails to complete the task due to strict safety alignment), whereas fine-tuning brings LLMs to a level comparable with encoder-based models. In terms of efficiency, encoder models finish inference within 20 seconds on both datasets, while LLMs require 2–9 minutes, reflecting differences in model scale, architecture, and the use of LoRA adapters.

4.3 Applicability to LLMs

Since LLMs are also Transformer-based, we conduct an exploratory study on transferring ToxiTrace to decoder-only LLMs. Due to resource constraints, LoRA was applied for parameter-efficient adaptation; we then replaced the decoding head with a binary classification head and optimized the same ToxiTrace objectives. As shown

Models	R	P	F_1	Inference Time
Llama3.2-3B	41.03	71.02	52.01	08m 06s
Qwen3-0.6B	73.68	67.10	70.23	09m 20s
Qwen3-1.7B	77.53	<u>71.48</u>	74.38	09m 28s
Llama-3.1-8B	81.32	69.21	74.78	12m 56s
Qwen3-8B	<u>84.83</u>	71.97	<u>77.87</u>	14m 33s
RoBERTa	42.37	68.42	52.34	01m 58s
RoBERTa*	71.27	59.87	65.08	01m 58s
RoBERTa + ours*	86.36	70.95	77.90	01m 58s

Table 2: Toxic span extraction performance on CNTP. Results without * correspond to extracting the top-15% most salient tokens, while results with * use our proposed bidirectional algorithm to extract high-saliency spans.

in Table 1, ToxiTrace achieves performance comparable to instruction fine-tuning on LLM, yet still falls short of encoder-based models trained with ToxiTrace. One possible reason is that LoRA updates only a tiny fraction of parameters (typically $\leq 0.5\%$), limiting its ability to reshape embedding-level gradients distributions required by our gradient-oriented training. We leave a systematic study of more effective parameter-efficient gradient shaping for future work.

4.4 Toxic Span Extraction Performance

We evaluate toxic span extraction on CNTP, where a prediction is counted as correct if its overlap with the gold span exceeds 50% (e.g., ground-truth toxic span: "河南人", pred: "河南" is correct, while "南" is incorrect). As shown in Table 2, replacing top-15% token selection with BiCSE yields a large gain for RoBERTa (RoBERTa \rightarrow RoBERTa*), confirming the benefit of extracting *contiguous* spans. Training with ToxiTrace further improves RoBERTa* to RoBERTa+ours* achieves an additional 12.82% F_1 -score improvement, indicating that CuSA and gradient-shaping objectives (GCLOSS/ARCL) produce more evidence-aligned saliency for extraction. Across LLMs, extraction generally improves with model size, but RoBERTa+ours* achieves comparable or better F_1 than the best LLM (Qwen3-8B) while requiring only about 1/7 of its inference time, demonstrating a substantially better accuracy–efficiency trade-off, the detailed prompt template is provided in Appendix C.

4.5 Ablation Study

We investigate the contributions of the key components in ToxiTrace through an ablation study with three variants, as shown in Table 3. Specifically,

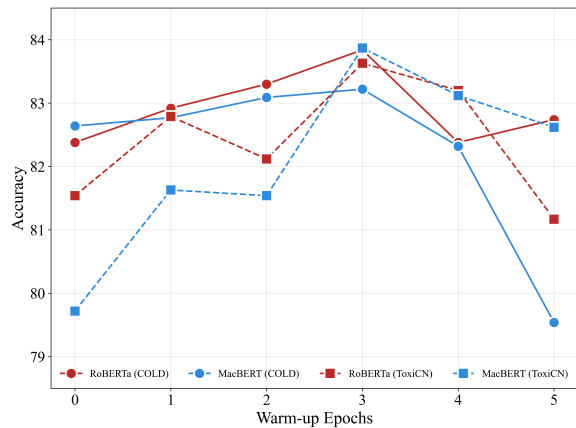


Figure 3: Final Accuracy with different Warm-up Epochs. The models achieve optimal classification performance when the warm-up steps are set to 3.

Full denotes the complete model that jointly optimizes GCLOSS and ARCL, while *w/o ARCL* and *w/o GCLOSS* remove the adversarial reasoning contrastive objective and the gradient constraint loss, respectively. We further include the vanilla *RoBERTa* as the baseline.

Since the ablation trends are consistent across COLL and ToxiCN, we only report results on COLL due to space constraints. Overall, removing either module leads to performance degradation in both classification and extraction. For classification, removing ARCL reduces Macro- F_1 by 0.56% (and slightly lowers F_1), indicating that ARCL provides a consistent gain via semantic regularization. Removing GCLOSS yields a smaller but noticeable drop in Macro- F_1 (0.32%), while increasing recall and decreasing precision, suggesting that without gradient shaping the model becomes more prone to over-predict toxic instances and thus sacrifices precision.

For toxic span extraction, GCLOSS has a substantially larger impact than ARCL. Removing ARCL causes a moderate degradation (extraction F_1 drops by 2.74%, with decreases in both recall and precision), whereas removing GCLOSS leads to a sharp decline (extraction F_1 drops by 12.75%, accompanied by large drops in recall and precision). Finally, compared with the RoBERTa baseline, the full model improves classification Macro- F_1 by 0.92% and boosts extraction F_1 by 12.82%, demonstrating the effectiveness of our joint training strategy.

Models	Classification					Extraction		
	Acc	R	P	F ₁	Macro-F ₁	R	P	F ₁
Full	83.84	86.19	76.14	80.85	83.68	86.36	70.95	77.90
w/o ARCL	83.13	86.71	74.72	80.27	83.12	82.55	68.99	75.16
w/o GCLoss	83.20	88.04	74.29	80.58	83.36	75.89	57.07	65.15
RoBERTa	82.75	86.43	74.24	79.87	82.76	71.27	59.87	65.08

Table 3: Ablation study results. "Classification" denotes results for the binary classification task on the COLD dataset; "Extraction" denotes results using the toxic span annotations in CNTP (Yang et al., 2025) as ground truth labels.

4.6 Effect of Warm-up Steps

Since our training pipeline requires warming up the foundation model for several epochs, this section analyzes the impact of the number of warm-up steps on the final classification and extraction results.

The number of warm-up steps determines the extent to which the model learns toxicity classification solely from binary labels. As shown in Figure 3, both too few and too many warm-up steps lead to a decrease in the final classification accuracy.

4.7 Gradient Attribution Faithfulness

In NLP, the faithfulness of explanations concerns whether the highlighted evidence truly reflects the causal decision process, rather than merely correlating with its prediction. Recent work argues that faithful explanations should be grounded in counterfactual reasoning and formalizes this intuition via *order-faithfulness*, showing that non-causal feature-importance methods can mis-rank evidence under confounding effects (Gat et al., 2024). Motivated by this view, we evaluate our gradient-based toxic span explanations from a causal perspective: if the extracted spans constitute genuine decision evidence, masking them should substantially reduce the model’s predicted toxicity confidence.

We extract toxic spans on the COLD dataset using two RoBERTa models, both localized by BiCSE: (1) a baseline RoBERTa trained with binary cross-entropy loss, and (2) a RoBERTa trained with our full objective in Eq. (9).

Figures 4 and 5 report the confidence drop, defined as the decrease in the predicted toxicity probability after masking the extracted spans. Compared with the BCE-only baseline, the ToxiTrace-trained model exhibits a consistently larger confidence reduction when its predicted spans are masked, indicating that the extracted spans are

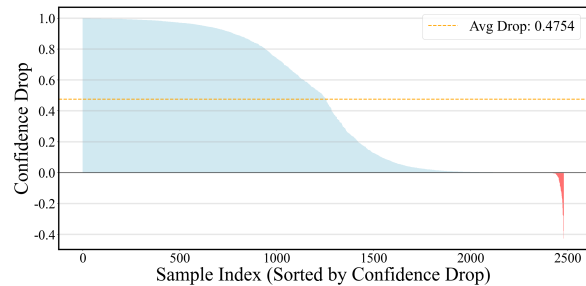


Figure 4: Confidence drop after masking BiCSE-extracted toxic spans for the RoBERTa baseline.

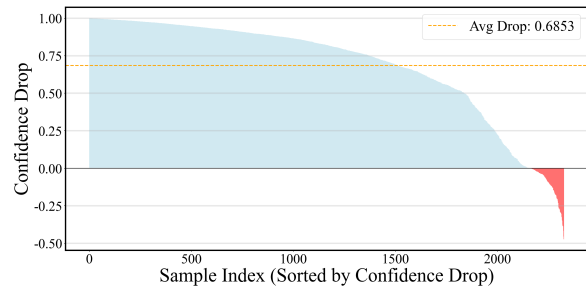


Figure 5: Confidence drop after masking BiCSE-extracted toxic spans for RoBERTa trained with ToxiTrace.

more decision-critical and thus more faithful to the model’s prediction.

5 Conclusion

We proposed ToxiTrace, which can automatically produce span-level annotations when fine-grained Chinese toxicity labels are unavailable. It further introduces a gradient-based loss to increase the model’s gradient responses on fine-grained toxic spans, and aligns each sentence with its corresponding reasoning in the semantic space, achieving simultaneous improvements in classification accuracy and interpretability. In addition, we design an algorithm for extracting high-saliency tokens, which addresses the limitation of prior methods that cannot recover contiguous high-saliency spans.

577 Limitations

578 In real-world scenarios, some toxic speeches are
579 "cloaked" or obfuscated, such as through the
580 use of homophones or Pinyin replacements (as
581 noted in CNTP). Most models experience a perfor-
582 mance degradation when processing such obfus-
583 cated toxic information. While datasets for this
584 exist in the Chinese domain, our current method
585 does not specifically address these perturbations.
586 We intend to incorporate robustness against such
587 variations into the scope of our future research.

588 Our method has been developed and evalu-
589 ated only on Chinese toxic content, which has
590 unique linguistic properties (e.g., character-level
591 tokenization and ambiguous semantic units). As
592 a result, its effectiveness on languages with differ-
593 ent structures remains to be verified.

594 Ethical Considerations

595 This work addresses toxic content detection and
596 explanation, which necessarily involves exposure
597 to offensive and harmful language. Although the
598 proposed method is designed to improve the faith-
599 fulness and transparency of model explanations,
600 there is a potential risk that fine-grained toxic span
601 extraction could be misused to reverse-engineer
602 moderation systems or to craft adversarial toxic
603 expressions that evade detection. To mitigate this
604 risk, we emphasize that the proposed framework
605 is intended for research and system auditing pur-
606 poses, rather than as a fully automated content
607 moderation solution, and should be deployed with
608 appropriate human oversight.

609 All datasets used in this study (COLD, Tox-
610 iCN, and CNTP) are publicly available benchmark
611 datasets released for academic research. They
612 do not contain personally identifying information,
613 and no additional data collection is conducted in
614 this work. While the datasets do include offensive
615 and toxic language by design, they are used solely
616 for model training and evaluation in a controlled
617 research setting. We do not attempt to identify, tar-
618 get, or profile any individuals, and no user-level or
619 sensitive attributes are inferred or stored.

620 References

621 Arnav Arora, Preslav Nakov, Momchil Hardalov,
622 Sheikh Muhammad Sarwar, Vibha Nayak, Yoan
623 Dinkov, Dimitrina Zlatkova, Kyle Dent, Ameya
624 Bhatawdekar, Guillaume Bouchard, and 1 others.

2023. Detecting harmful content on online plat-
forms: what platforms need vs. where research ef-
forts go. *ACM Computing Surveys*, 56(3):1–17. 625
626
627

Sylvia W Azumah, Nelly Elsayed, Zag ElSayed, and
Murat Ozer. 2023. Cyberbullying in text content de-
tection: an analytical review. *International Journal
of Computers and Applications*, 45(9):579–586. 628
629
630
631

Pinkesh Badjatiya, Shashank Gupta, Manish Gupta,
and Vasudeva Varma. 2017. Deep learning for hate
speech detection in tweets. In *Proceedings of the
26th International Conference on World Wide Web
Companion*, WWW '17 Companion, page 759–760. 632
633
634
635
636

Jasmijn Bastings, Sebastian Ebert, Polina Zablotskaia,
Anders Sandholm, and Katja Filippova. 2022. "will
you find these shortcuts?" a protocol for evaluating
the faithfulness of input salience methods for text
classification. In *Proceedings of the 2022 Confer-
ence on Empirical Methods in Natural Language
Processing*, pages 976–991. 637
638
639
640
641
642
643

Tommaso Caselli, Valerio Basile, Jelena Mitrović, and
Michael Granitzer. 2021. HateBERT: Retraining
BERT for abusive language detection in English. In
*Proceedings of the 5th Workshop on Online Abuse
and Harms (WOAH 2021)*, pages 17–25. 644
645
646
647
648

August F.Y. Chao, Chen-Shu Wang, Bo-Yi Li, and
Hong-Yan Chen. 2024. From hate to harmony:
Leveraging large language models for safer speech
in times of covid-19 crisis. *Heliyon*, 10(16):e35468. 649
650
651
652

George Chrysostomou and Nikolaos Aletras. 2021a.
Enjoy the salience: Towards better transformer-
based faithful explanations with word salience. In
*Proceedings of the 2021 Conference on Empirical
Methods in Natural Language Processing*, pages
8189–8200. 653
654
655
656
657
658

George Chrysostomou and Nikolaos Aletras. 2021b.
Improving the faithfulness of attention-based expla-
nations with task-specific information for text clas-
sification. In *Proceedings of the 59th Annual Meet-
ing of the Association for Computational Linguistics
and the 11th International Joint Conference on Nat-
ural Language Processing (Volume 1: Long Papers)*,
pages 477–488. 659
660
661
662
663
664
665
666

Yu-Neng Chuang, Guanchu Wang, Chia-Yuan Chang,
Ruixiang Tang, Shaochen Zhong, Fan Yang, Meng-
nan Du, Xuanting Cai, Vladimir Braverman, and
Xia Hu. 2025. Faithlm: Towards faithful ex-
planations for large language models. *Preprint*,
arXiv:2402.04678. 667
668
669
670
671
672

Antonia Creswell and Murray Shanahan. 2022. Faith-
ful reasoning using large language models. *Preprint*,
arXiv:2208.14271. 673
674
675

Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, Shi-
jin Wang, and Guoping Hu. 2020. Revisiting pre-
trained models for Chinese natural language pro-
cessing. In *Findings of the Association for Compu-
tational Linguistics: EMNLP 2020*, pages 657–668. 676
677
678
679
680

681	Thomas Davidson, Dana Warmley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language . <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , 11(1):512–515.	736
682		737
683		738
684		739
685		
686	Jiawen Deng, Zhuang Chen, Hao Sun, Zhexin Zhang, Jincenzi Wu, Satoshi Nakagawa, Fuji Ren, and Minlie Huang. 2023. Enhancing offensive language detection with data augmentation and knowledge distillation . <i>Research</i> , 6:0189.	740
687		741
688		742
689		743
690		744
691	Jiawen Deng, Jingyan Zhou, Hao Sun, Chujie Zheng, Fei Mi, Helen Meng, and Minlie Huang. 2022. COLD: A benchmark for Chinese offensive language detection . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11580–11599.	745
692		746
693		747
694		748
695		749
696		750
697	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186.	751
698		752
699		753
700		754
701		755
702		
703		
704		
705	Shuoyang Ding and Philipp Koehn. 2021. Evaluating saliency methods for neural language models . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5034–5052.	756
706		757
707		758
708		759
709		
710		
711	Tianqing Fang, Wenxuan Zhou, Fangyu Liu, Hongming Zhang, Yangqiu Song, and Muhao Chen. 2024. On-the-fly denoising for data augmentation in natural language understanding . In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 766–781.	760
712		761
713		762
714		763
715		764
716		
717	Yair Gat, Nitay Calderon, Amir Feder, Alexander Chapanin, Amit Sharma, and Roi Reichart. 2024. Faithful explanations of black-box nlp models using llm-generated counterfactuals . In <i>International Conference on Representation Learning</i> , volume 2024, pages 48946–48979.	765
718		766
719		767
720		768
721		769
722		770
723	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models . <i>ICLR</i> , 1(2):3.	771
724		772
725		773
726		774
727	Xuming Hu, Yong Jiang, Aiwei Liu, Zhongqiang Huang, Pengjun Xie, Fei Huang, Lijie Wen, and Philip S. Yu. 2023. Entity-to-text based data augmentation for various named entity recognition tasks . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 9072–9087.	775
728		776
729		777
730		778
731		779
732		780
733	Sérgio Jesus, Catarina Belém, Vladimir Balayan, João Bento, Pedro Saleiro, Pedro Bizarro, and João Gama. 2021. How can i choose an explainer? an application-grounded evaluation of post-hoc explanations . In <i>Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21</i> , page 805–815.	781
734		782
735		
	Jonathan Kamp, Lisa Beinborn, and Antske Fokkens. 2023. Dynamic top-k estimation consolidates disagreement between feature attribution methods . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6190–6197.	783
		784
		785
		786
		787
		788
	Hannah Kirk, Abeba Birhane, Bertie Vidgen, and Leon Derczynski. 2022. Handling and presenting harmful text in NLP research . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 497–510.	789
		790
	Satyapriya Krishna, Tessa Han, Alex Gu, Steven Wu, Shahin Jabbari, and Himabindu Lakkaraju. 2025. The disagreement problem in explainable machine learning: A practitioner’s perspective . <i>Preprint</i> , arXiv:2202.01602.	
	Irene Kwok and Yuzhou Wang. 2013. Locate the hate: Detecting tweets against blacks . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 27(1):1621–1622.	
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach . <i>arXiv preprint arXiv:1907.11692</i> .	
	Junyu Lu, Bo Xu, Xiaokun Zhang, Changrong Min, Liang Yang, and Hongfei Lin. 2023. Facilitating fine-grained detection of Chinese toxic language: Hierarchical taxonomy, resources, and benchmarks . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16235–16250.	
	Behrad Moniri, Hamed Hassani, and Edgar Dobriban. 2025. Evaluating the performance of large language models via debates . In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 2040–2075.	
	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?": Explaining the predictions of any classifier . In <i>Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16</i> , page 1135–1144.	
	Andrew Slavin Ross, Michael C. Hughes, and Finale Doshi-Velez. 2017. Right for the right reasons: training differentiable models by constraining their explanations . In <i>Proceedings of the 26th International Joint Conference on Artificial Intelligence, IJCAI '17</i> , page 2662–2670.	
	Evgenia Rusak, Patrik Reizinger, Attila Juhas, Oliver Bringmann, Roland S. Zimmermann, and Wieland	

791	Brendel. 2025. Infonce: Identifying the gap between theory and practice . In <i>Proceedings of The 28th International Conference on Artificial Intelligence and Statistics</i> , volume 258 of <i>Proceedings of Machine Learning Research</i> , pages 4159–4167.	848
792		849
793		850
794		851
795		852
796	Diptanu Sarkar, Marcos Zampieri, Tharindu Ranasinghe, and Alexander Ororbia. 2021. fBERT: A neural transformer for identifying offensive content . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1792–1798.	853
797		854
798		855
799		856
800		857
801	Maximilian Schmidhuber and Udo Kruschwitz. 2024. LLM-based synthetic datasets: Applications and limitations in toxicity detection . In <i>Proceedings of the Fourth Workshop on Threat, Aggression & Cyberbullying @ LREC-COLING-2024</i> , pages 37–51.	858
802		859
803		860
804		861
805		862
806	Sandipan Sikdar, Parantapa Bhattacharya, and Kieran Heese. 2021. Integrated directional gradients: Feature interaction attribution for neural NLP models . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 865–878.	863
807		864
808		865
809		866
810		867
811		868
812		869
813		870
814	Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shangwei Guo, Tianwei Zhang, and Guoyin Wang. 2023. Text classification via large language models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 8990–9005.	871
815		872
816		873
817		874
818		875
819	Zijun Sun, Xiaoya Li, Xiaofei Sun, Yuxian Meng, Xiang Ao, Qing He, Fei Wu, and Jiwei Li. 2021. ChineseBERT: Chinese pretraining enhanced by glyph and Pinyin information . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2065–2075.	876
820		877
821		878
822		
823		
824		
825		
826		
827	Gemini Team. 2025. Gemini: A family of highly capable multimodal models . <i>Preprint</i> , arXiv:2312.11805.	
828		
829		
830	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	
831		
832		
833		
834		
835	Zeeraq Waseem and Dirk Hovy. 2016. Hateful symbols or hateful people? predictive features for hate speech detection on Twitter . In <i>Proceedings of the NAACL Student Research Workshop</i> , pages 88–93.	
836		
837		
838		
839	Tomer Wullach, Amir Adler, and Einat Minkov. 2022. Character-level hypernetworks for hate speech detection . <i>Expert Systems with Applications</i> , 205:117571.	
840		
841		
842		
843	Qiyao Xue, Yuchen Dou, Ryan Shi, Xiang Lorraine Li, and Wei Gao. 2025. Mmbert: Scaled mixture-of-experts multimodal bert for robust chinese hate speech detection under cloaking perturbations . <i>Preprint</i> , arXiv:2508.00760.	
844		
845		
846		
847		
	Shujian Yang, Shiyao Cui, Chuanrui Hu, Haicheng Wang, Tianwei Zhang, Minlie Huang, Jialiang Lu, and Han Qiu. 2025. Exploring multimodal challenges in toxic Chinese detection: Taxonomy, benchmark, and findings . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 14382–14396.	
	Bing Zhang, Mikio Takeuchi, Ryo Kawahara, Shubhi Asthana, Md. Maruf Hossain, Guang-Jie Ren, Kate Soule, Yifan Mai, and Yada Zhu. 2025. Evaluating large language models with enterprise benchmarks . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 3: Industry Track)</i> , pages 485–505.	
	Yaosheng Zhang, Tiegang Zhong, Tingjun Yi, and Haoming Li. 2024. Domain-enhanced prompt learning for chinese implicit hate speech detection . <i>IEEE Access</i> , 12:13773–13782.	
	Kun Zhou, Wayne Xin Zhao, Sirui Wang, Fuzheng Zhang, Wei Wu, and Ji-Rong Wen. 2021. Virtual data augmentation: A robust and general framework for fine-tuning pre-trained models . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 3875–3887.	
	Steven Zimmerman, Udo Kruschwitz, and Chris Fox. 2018. Improving hate speech detection with deep learning ensembles . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> .	

879 A Algorithm

880 The algorithm pseudo-code to extract continuous
881 token spans.

Algorithm 1 Bidirectional Salient Span Extraction

Input: Gradient sequence $G = \{g_1, g_2, \dots, g_n\}$

Output: Salient span set S

```

1:  $\mu \leftarrow \text{Mean}(G)$ 
2:  $\tau \leftarrow \text{Median}(|g_i - g_{i-1}|)$  for  $i \in [2, n]$ 
3:
4: Function FindCliffEnd( $G, \text{start}, \mu, \tau$ ):
5:    $p \leftarrow \text{start}$ 
6:   for  $i = \text{start}$  to  $n$  do
7:     if  $g_i > \mu$  then  $p \leftarrow i$ 
8:     if  $i \leq n - 2$  and  $g_i - g_{i+1} > \tau$  then
9:       if  $g_{i+1} - g_{i+2} \leq \tau$  then return  $i$ 
10:    if  $g_i \leq \mu$  then return  $p$ 
11:  return  $p$ 
12:
13: Function ForwardScan( $G, \mu, \tau$ ):
14:   $S \leftarrow \emptyset, i \leftarrow 2$ 
15:  while  $i \leq n$  do
16:    if  $g_i > \mu$  and  $g_i - g_{i-1} > \tau$  then
17:       $s \leftarrow i, e \leftarrow \text{FindCliffEnd}(G, s, \mu, \tau)$ 
18:       $S \leftarrow S \cup \{(s, e)\}, i \leftarrow e + 1$ 
19:    else
20:       $i \leftarrow i + 1$ 
21:  return  $S$ 
22:
23:  $S_{\text{fwd}} \leftarrow \text{ForwardScan}(G, \mu, \tau)$ 
24:  $G \leftarrow \text{Reverse}(G)$ 
25:  $S_{\text{bwd}} \leftarrow \text{Reverse}(\text{ForwardScan}(G, \mu, \tau))$ 
26: return  $\text{Merge}(S_{\text{fwd}} \cup S_{\text{bwd}})$ 

```

882 **Threshold Computation** Given a gradient se-
883 quence $G = g_1, g_2, \dots, g_n$, we compute two
884 thresholds:

885 Mean threshold $\mu = \text{Mean}(G)$, serving as the
886 baseline for identifying high-gradient tokens. Dif-
887 ference threshold $\tau = \text{Median}(|g_i - g_{i-1}|)$, captur-
888 ing the typical magnitude of gradient transitions.

889 **Start Condition** A span begins at position i if
890 the gradient exhibits a steep ascent:

$$g_i > \mu \quad \text{and} \quad g_i - g_{i-1} > \tau$$

891 **Termination Condition 1 (Cliff Edge Detection)**

892 The span terminates at position i when a cliff
893 edge is detected—i.e., the current position shows

a steep descent but the subsequent descent dimin-
ishes:

$$g_i - g_{i+1} > \tau \quad \text{and} \quad g_{i+1} - g_{i+2} \leq \tau$$

This condition identifies the boundary where
the gradient "cliff" ends, ensuring complete span
extraction by continuing to search until the final
cliff edge is found.

Termination Condition 2 (Fallback) If no cliff
edge is detected, the span ends at the last position
where $g_i > \mu$, preventing incomplete extraction
when gradients decay gradually.

Bidirectional Scanning To capture spans that
may be missed by unidirectional scanning, we per-
form forward scanning on G and backward scan-
ning on the reversed sequence G' . The final span
set is obtained by merging overlapping intervals
from both directions:

$$S = \text{Merge}(S_{\text{fwd}} \cup S_{\text{bwd}})$$

This bidirectional approach ensures robust de-
tection of salient spans regardless of their orienta-
tion within the sequence.

B Changes in Saliency Map

913

Models trained with our proposed method exhibit a pronounced shift in saliency maps, characterized by substantially higher saliency scores assigned to toxic spans. Figure 6 and 7 presents an example illustrating the saliency shift before and after training with proposed ToxiTrace, demonstrating that the model becomes more focused on truly toxic spans.

914

915

916

917

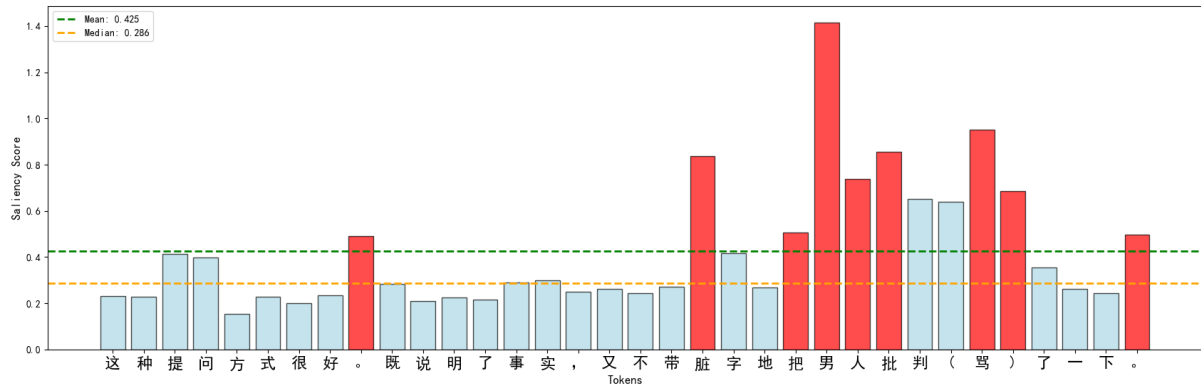


Figure 6: Saliency map of a RoBERTa model trained only with binary classification labels.

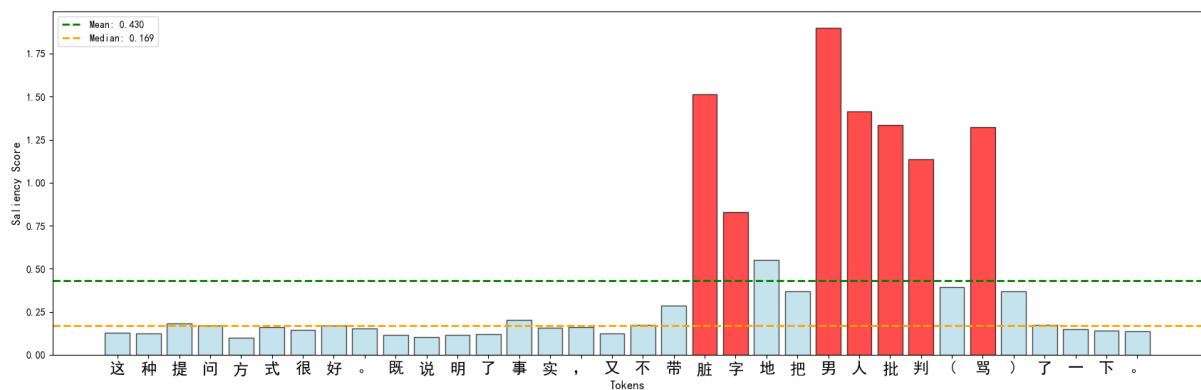


Figure 7: Saliency map of a RoBERTa model after training with our proposed method.

918
919
920
921
922

C Prompt Templates for LLMs

In this study, LLMs are employed for multiple tasks. Table 4 presents the prompt templates used for toxic span extraction. Tables 7 and 9 provide the prompts designed to elicit sample-specific adversarial explanations for toxic and non-toxic instances, respectively. Table 11 summarizes the prompts used to prompt LLMs to directly assess the toxicity of input texts.

Toxic Spans Refining Prompt

System: 你是一个有害内容识别领域的专家，在你目前对“有害内容”的定义的基础之上，请对以下这个用户提供的内容进行有害内容实体的片段定位任务，实体片段尽量短且尽量包含多的有害信息（请参考后面的输出范例，所有输入句都**一定存在这样的片段**，每句话都会给出句中可疑的一些字段作为线索供参考）。
片段之间不要有重合或包含，也要注意中文隐式有害俚语。输出不需要给出任何解释、不要进行任何额外说明。输出示例（使用单行形式，不要套入 markdown 格式）：
有害实体 1、有害实体 2、有害实体 3

真实数据示例 1:
用户输入：确实现在的拳师天生就是坏种它们贪婪傲慢懒惰暴食淫欲嫉妒暴怒七宗罪它们全部都犯了它们欲壑难填利令智昏是畜生不如的东西
输出：拳师、是坏种、淫欲、七宗罪、都犯了、欲壑难填、畜生不如

真实数据示例 2:
用户输入：厦门河南人和江西人开出租的很多，都挺狠的
输出：河南人、江西人、挺狠

真实数据示例 3:
用户输入：这帮黑人不仅有枪还有毒品，一贫如洗？
输出：黑人、有枪、有毒品

User:
{sentence}
【参考线索（可能相关的重点词汇）】: {hint_text}

Table 4: Toxic Spans Refining Prompt

Translation of Toxic Spans Refining Prompt

System:
You are an expert in the field of toxic content recognition. Based on your current definition of "toxic content," please perform a toxic entity span extraction task on the content provided by the user below. The entity spans should be as short as possible while containing as much toxic information as possible.

(Please refer to the output examples below. All input sentences **definitely contain such fragments**, and some suspicious fields in the sentence will be given as clues for reference).

There should be no overlap or containment between fragments. Pay attention to implicit toxic slang.
Output Requirement: Do not provide any explanation or extra description.
Output Example (use a single line format, do not wrap in markdown): Toxic Entity 1, Toxic Entity 2, Toxic Entity 3

Real Data Example 1:
User Input: Indeed, current feminists are born bad seeds; they are greedy, arrogant, lazy, gluttonous, lustful, jealous, and wrathful—committing all seven deadly sins. Their desires are insatiable, they are blinded by greed, and they are less than beasts.
Output: feminists, are bad seeds, lustful, seven deadly sins, committing all, desires are insatiable, less than beasts

Real Data Example 2:
User Input: There are many taxi drivers from Henan and Jiangxi in Xiamen, and they are all pretty ruthless.
Output: Henan people, Jiangxi people, pretty ruthless

Real Data Example 3:
User Input: These black people not only have guns but also drugs, and yet they claim to be destitute?
Output: black people, have guns, have drugs

User:
{sentence}
【Reference Clues (Potential Key Vocabulary)】: {hint_text}

Table 5: Translation of Toxic Spans Refining Prompt

Translation of Toxic Spans Refining Prompt

System:

You are an expert in the field of toxic content recognition. Based on your current definition of "toxic content," please perform a toxic entity span extraction task on the content provided by the user below. The entity spans should be as short as possible while containing as much toxic information as possible.

(Please refer to the output examples below. All input sentences **definitely contain such fragments**, and some suspicious fields in the sentence will be given as clues for reference).

There should be no overlap or containment between fragments. Pay attention to implicit toxic slang.

Output Requirement: Do not provide any explanation or extra description.

Output Example (use a single line format, do not wrap in markdown): Toxic Entity 1, Toxic Entity 2, Toxic Entity 3

Real Data Example 1:

User Input: Indeed, current feminists are born bad seeds; they are greedy, arrogant, lazy, gluttonous, lustful, jealous, and wrathful—committing all seven deadly sins. Their desires are insatiable, they are blinded by greed, and they are less than beasts.

Output: feminists, are bad seeds, lustful, seven deadly sins, committing all, desires are insatiable, less than beasts

Real Data Example 2:

User Input: There are many taxi drivers from Henan and Jiangxi in Xiamen, and they are all pretty ruthless.

Output: Henan people, Jiangxi people, pretty ruthless

Real Data Example 3:

User Input: These black people not only have guns but also drugs, and yet they claim to be destitute?

Output: black people, have guns, have drugs

User:

{sentence}

【Reference Clues (Potential Key Vocabulary)】 : {hint_text}

Table 6: Translation of Toxic Spans Refining Prompt

Toxic Reasoning Generation Prompt

你是一个中文毒性内容理解的专家，用户将给定文本，请完全站在“这段内容有害”的主观立场上认定该内容包含对他人的攻击/威胁/贬损/歧视等“有害表达”。

请只输出一句简洁理由，不要有任何多余的内容，需要针对这句话的内容进行针对性的解释，不要输出宽泛、模糊的解释。

参考输出（仅作为示例，并非输出语句格式模版，但需至少包含其中提到针对该句的要素，即群体、观点、用词、可能的负面理解等）：

这段话包含对[某一群体]的[xx]性质的贬损性描述，使用了“[xxx]”这样的词汇，容易被理解为嘲讽或不尊重，也可能被理解为人身攻击。

输入：{sentence}

Table 7: Toxic Reasoning Generation Prompt

Translation of Toxic Reasoning Generation Prompt

You are an expert in understanding toxic content. The user will provide text. Please stand completely on the subjective standpoint that "this content is toxic" and determine that the content contains "toxic expressions" such as attacks, threats, derogation, or discrimination against others.

Please output only one concise reason. Do not include any superfluous content. You need to provide a specific explanation targeting the content of this sentence; do not output broad or vague explanations.

Reference Output (This is **not** a template, but the output must include elements targeting the sentence, such as the group, viewpoint, wording, and possible negative interpretations):

This statement contains a derogatory description of the [xx] nature regarding [specific group], using vocabulary like "[xxx]", which is easily understood as mockery or disrespect, and may also be interpreted as a personal attack.

Input: {sentence}

Table 8: Translation of Toxic Reasoning Generation Prompt

Non-toxic Reasoning Generation Prompt

用户将给定文本，请完全站在“这段内容完全正常”的主观立场上认定该内容完全不包含对他人的攻击/威胁/贬损/歧视等任何“有害表达”。

请只输出一句简洁理由，不要有任何多余的内容，需要针对这句话的内容进行针对性的解释，不要输出宽泛、模糊的解释。

参考输出（并非输出语句格式模版，但需包含其中提到针对该句的要素，即群体、观点、用词、不可能产生负面理解的原因等）：

这段话针对 [xx] 群体发表了 [客观观点]，表达一种情绪宣泄，即使使用了 “[xxx]” 这样的词汇，也只是语气比较重，没有贬损的意味。

输入：{sentence}

Table 9: Non-toxic Reasoning Generation Prompt

Translation of Toxic Spans Refining Prompt

The user will provide text. Please stand completely on the subjective standpoint that "this content is completely normal" and determine that the content does not contain any "toxic expressions" such as attacks, threats, derogation, or discrimination against others.

Please output only one concise reason. Do not include any superfluous content. You need to provide a specific explanation targeting the content of this sentence; do not output broad or vague explanations.

Reference Output (This is **not** a template, but the output must include elements targeting the sentence, such as the group, viewpoint, wording, and reasons why it cannot generate negative interpretations):

This statement expresses an [objective viewpoint] regarding the [xx] group, serving as an emotional outlet. Even though it uses vocabulary like "[xxx]", it is merely a strong tone and carries no derogatory implication.

Input: {sentence}

Table 10: Translation of Toxic Reasoning Generation Prompt

LLM Classification Task Instruction

你是一个有害内容识别领域的专家，现在有一些潜在的有害内容文本需要你来分析，可能的有害内容包括但不限于：

地域歧视：针对某些地区或该地区的人发表的整体贬低性言论；

性别对立：表达针对某一性别的贬低或歧视性言论；

种族歧视：针对特定人种发表的歧视性表达。

其余未列在上述分类中的如果有明显的攻击性及仇恨言论，也认为是有害内容。

你只需要对用户提供的内容进行判断，然后对你认为有害的内容输出一个“1”，无害内容输出“0”即可，无需进行任何其他分析，只需要输出一个“0”或“1”即可。

输入：{sentence}

Table 11: LLM Classification Task Instruction

Translation of LLM Classification Task Instruction

You are an expert in the field of toxic content recognition. There are some potentially toxic texts that need your analysis. Possible toxic content includes but is not limited to:

Regional Discrimination: Overall derogatory remarks against specific regions or people from those regions; Gender Opposition: Expressing derogatory or discriminatory remarks against a specific gender; Racial Discrimination: Discriminatory expressions against specific races.

If there are obvious offensive or hate speech remarks not listed in the above categories, they are also considered toxic content.

You only need to judge the content provided by the user. Output a "1" for content you consider toxic, and "0" for non-toxic content. Any other analysis is prohibited; strictly output only a "0" or "1".

Input: {sentence}

Table 12: Translation of LLM Classification Task Instruction