

HARM: Learning Hate-Aware Reward Model for Evaluating Natural Language Explanations of Offensive Content

Lorenzo Puppi Vecchi¹, Alceu de Souza Britto Jr¹,
Emerson Cabrera Paraiso¹, Rafael Menelau Cruz²,

¹Graduate Program in Informatics - Pontifícia Universidade Católica do Paraná, Brasil,

²Département de Génie Logiciel et TI - École de Technologie Supérieure ÉTS, Canada

Abstract

Explaining why content is hateful using natural language is crucial for fostering transparency in automated content moderation systems. However, evaluating the quality of such explanations remains an open challenge. General-purpose reward models (RMs), commonly used for scoring natural language outputs, are typically optimized for broad notions of safety. We argue that this optimization penalizes situations where references to stereotypes or offensive content are essential for explanations with higher explanatory fidelity. To address this gap, we introduce *SBIC-Explain*, a human-validated dataset of 370,788 LLM generated NLEs for offensive content, spanning three levels of human-annotated contextual richness: *Tier 1*: text-only, *Tier 2*: + classification-aware, and *Tier 3*: + semantics-informed. We hypothesize that as human-annotated context increases, explanations should lead to higher perceived explanations with higher explanatory fidelity. Yet, we find that existing RMs systematically assign lower scores to more contextually rich (and often more offensive) explanations, revealing a misalignment between model preferences and explanatory fidelity for this context. We propose *HARM* (*Hate-Aware Reward Model*), a RM that integrates interpretable signals to better align reward scores with the needs of hate speech explanation. HARM outperforms general-purpose baselines, improving NLE pair-wise preference. Available at: <https://github.com/Lorenzo815/HARM>.

1 Introduction

The proliferation of hate speech on social media poses substantial risks to individuals and communities, threatening social cohesion and even inciting real-world violence (Jahan and Oussalah, 2021; Malik et al., 2023). In response, automated hate speech detection systems have become essential for moderating online content and ensuring safer digital environments (Narula and Chaudhary, 2024).

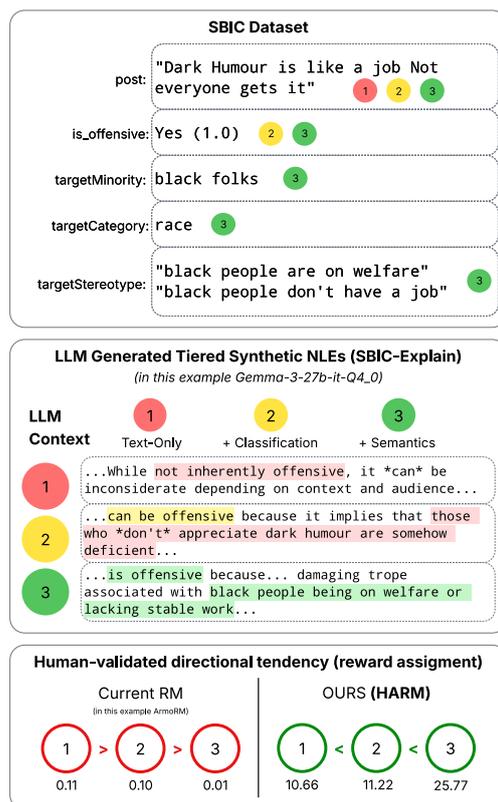


Figure 1: **Concept Overview.** We leverage a human annotated dataset (top) to generate LLM synthetic NLEs (middle) under three levels of contextual access. Intuitively, access to richer human context should yield explanations that receive higher reward scores. However, existing general-purpose reward models (RMs) assign lower scores as the language becomes more sensitive (bottom). Our proposed method, **HARM**, learns to better reflect the **human-validated directional tendency** that richer context is often preferred.

However, beyond mere classification, a central challenge lies in understanding why content is hateful. One promising solution is the use of AI-generated NLEs, which aim to make model reasoning interpretable, enhancing transparency and fostering user trust (Mathew et al., 2022).

Yet, effectively training and evaluating such explanations remains difficult, in part due to limita-

tions in current explainable hate speech datasets. Most existing resources rely on shallow signals, such as highlighting offensive spans (Mathew et al., 2022; Arshad and Shahzad, 2024; Delbari et al., 2024; Hoang et al., 2023; Pavlopoulos et al., 2021; Ravikiran and Annamalai, 2021), which lack the depth needed to support rich, contextual, and stereotype-aware explanations (Sap et al., 2020). While useful for pinpointing relevant tokens, such approaches often fail to capture the broader social context or implicit stereotypes that underpin many harmful messages. For example, highlighting the word ginger does little to surface its role as a slur against redheads. Similarly, in the statement Dark Humour is like a job Not everyone gets it (from Figure 1), the surface tone masks a stereotype associating racial groups with unemployment.

The Social Bias Inference Corpus (SBIC) (Sap et al., 2020) addresses these shortcomings by including human-written implied statements that surface the stereotypical or biased assumptions underlying ostensibly innocuous content. While SBIC was originally constructed to capture offensiveness, it offers a rich annotation framework for studying *implicit hate speech* (the kind that relies on euphemism, ambiguity, or stereotypes rather than explicit slurs (Kim et al., 2023)).

Despite the growing use of LLM-generated NLEs in hate speech detection (Huang et al., 2024, 2023a), existing datasets such as HateCOT (Nghiem and Daumé Iii, 2024) overlook this contextual depth. They do not systematically evaluate how access to rich, stereotype-aware human annotations, like those in SBIC, affects explanation quality. This gap raises a crucial question: *How much does access to human-annotated contextual information improve the quality of NLEs in hate speech detection?*

Current limitation in explanation quality evaluation is not just a data issue, it also stems from how explanation outputs are scored. Beyond their role in reinforcement learning for fine-tuning language models (e.g., RLHF), reward models (RMs) are increasingly used as scoring functions to evaluate the quality of LLM-generated content, including NLEs. However, general-purpose RMs are poorly aligned with the goals of hate speech explanation (Christian et al., 2025). Trained primarily to promote safety, they tend to penalize outputs that reference stereotypes or offensive language, even when such content is necessary for a post-grounded explanation. This misalignment incentivizes the production

of sanitized but less well-grounded explanations, which sacrifice nuance and clarity for safety compliance. The result reflects a deeper alignment dilemma: while harmful language must be curbed, explanations with higher fidelity often require referencing uncomfortable truths (Chua et al., 2024; Lyu et al., 2024).

To investigate these phenomena, we introduce *SBIC-Explain*, a dataset of synthetic natural language explanations (NLEs) grounded in SBIC’s human-annotated stereotypes and validated through a human study on a sample. We show that state-of-the-art models systematically undervalue explanations with higher hate speech fidelity, particularly those referencing offensive stereotypes. To bridge this gap, we propose a lightweight method for building a domain-specific *Hate-Aware Reward Model (HARM)*, which re-weights interpretable outputs of general-purpose reward models to better evaluate hate-related content. HARM surpasses LLM-as-a-Judge based approaches as well as general purpose RMs, correctly ranking NLEs that have higher explanatory fidelity when explaining offensive content.

We organize our contributions as follows:

(i) We construct and release **SBIC-Explain**, a multi-tier, human validated, multi-model synthetic NLE dataset grounded in stereotype-level annotations, to benchmark the effect of contextual depth on hate speech explanation. (ii) We reveal a critical failure mode in current reward methods in the context of NLE for offensive content, showing how safety-oriented training biases them against truthful explanations. (iii) We propose **HARM**, a hate-aware reward model that uses a lightweight adapter to improve NLE evaluation for sensitive content.

2 Related Work

2.1 Explainable Hate Speech Datasets

In the landscape of existing explainable hate speech datasets, Table 5 in Appendix A reveal that most works treat explainability as a span annotation.

Predominance of Span Rationales. Seven out of ten listed datasets (HateBRXplain (Mathew et al., 2022), HateInsights (Arshad and Shahzad, 2024), PHate (Delbari et al., 2024), ViHOS (Hoang et al., 2023), HateXplain (Mathew et al., 2022), SemEval2021 (Pavlopoulos et al., 2021), DOSA (Ravikiran and Annamalai, 2021)) rely primarily on *span rationales*, where annotators highlight token-

level or phrase-level segments deemed offensive. While span rationales help models locate surface clues (e.g., explicit slurs), they frequently omit the broader sociohistorical context behind a slurs offensiveness. In practice, this limitation can lead to shallow explanations, a classifier may learn to flag a word without understanding why that term is harmful to a particular group.

Unique Role of SBICs FreeText Annotations. SBIC (Sap et al., 2020) stands out in Table 5 because it provides not only an offensiveness label (offensiveYN) and categorical fields (targetMinority, targetCategory) but also *freetext implied statements* capturing why a particular target is stigmatized. Although SBIC has been used for training classifiers and for categorization tasks, downstream explainability work has largely ignored the option it offers for systematically testing the impact of its 34K implied statements for explainability (Nghiem and Daumé Iii, 2024). This underutilization represents a missed opportunity: *freetext stereotypes offer semantic depth that neither span highlights nor simple labels can provide for grounding explanations.*

Therefore, this work aims to test how much the incorporation of such semantic depth improves the perceived quality of NLEs, particularly in the eyes of reward models. **Our hypothesis is that explanations grounded in richer, stereotype-informed context leads to higher perceived explanatory fidelity.**

2.2 Evaluating NLEs with Preference Modeling

Early approaches. With the aim to compare AI generated content with a human ground truth, early approaches relied on lexical overlap metrics like BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004), but these are poorly correlated with human judgments of explanation quality (Sai et al., 2020). More recent metrics based on pre-trained language models (PLMs), like BERTScore (Zhang et al., 2020), offer better semantic evaluation, but are notorious for carrying unfair stereotypical bias such as racial, gender, or religion bias (Kaneko and Bollegala, 2021; Sun et al., 2022). This behavior is unacceptable for a work that aims to specifically deal with hate speech.

Reward Models (RMs). Created to improve the challenge of aligning large language models

(LLMs) with human preferences, RMs were initially trained to predict a single scalar score for a given context, and generated continuation (Ziegler et al., 2020; Chen et al., 2024b). These approaches underpin alignment techniques such as Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022). However, human preferences are rarely monolithic; they are inherently multi-dimensional, encompassing trade-offs between attributes helpfulness and informativeness (Yu et al., 2025; Bai et al., 2022). A significant issue lies when dealing with sensitive or identity-related language (Sap et al., 2019), resulting in over censorship, linguistic erasure, and the suppression of valid discourse, even if necessary to convey intent with higher explanatory fidelity (Christian et al., 2025).

Multi-attribute RMs. As a natural evolution of single scalar RM, such approaches generate multiple scalars that disentangle genuine helpfulness from unimportant factors like length bias and offer precise, steerable rewards (Chen et al., 2024a). Datasets such as HelpSteer (Wang et al., 2023b) and UltraFeedback (Cui et al., 2024) provide fine-grained annotations across multiple dimensions. Building on these, ArmoRM (Wang et al., 2024a) learns to predict multi-objective reward scalars¹ for each response.

LLM-as-a-Judge (LaaJ). A parallel line of research to reward modeling is the use of *LLMs themselves as evaluators*, commonly termed *LLM-as-a-Judge (LaaJ)* (Kim, 2025). Instead of relying on fixed metrics, pretrained LLMs are prompted to act as judges, either through *absolute scoring* (evaluating one response against a rubric) or *pair-wise comparison* (choosing which of two responses better satisfies given criteria) (Li et al., 2024).

Recent advances have introduced more structured frameworks for LLM evaluation. The *G-Eval* framework (Liu et al., 2023) augments judgments with chain-of-thought rationales, improving both correlation with human ratings and interpretability. A growing line of work proposes dedicated judge models, such as *Prometheus* (Kim et al., 2024), which integrates reference answers and fine-grained rubrics, achieving strong agreement with human evaluators. More recently, *GPTScore* (Fu et al., 2024) enables instruction-guided and multi-dimensional evaluation without task-specific train-

¹Further discussed in Appendix J

ing, while *TIGERScore* (Jiang et al., 2024) leverages fine-tuning to yield a reference-free, high-correlation evaluator. Systems such as *Praetor* (Leng et al., 2025) further emphasize customizable, criterion-specific judgments, reflecting the trend toward flexible, rubric-grounded evaluators.

Despite these advances, persistent limitations remain. Studies highlight *position bias*, susceptibility to verbosity or formatting, and difficulties with calibration across domains (Wang et al., 2024b). Recent evaluations further suggest that fine-tuned judges often underperform frontier models like GPT-4 in cross-domain robustness (Huang et al., 2025).

Limitations of Automated Evaluation Approaches. To test both RMs and LaaJ, leading benchmarks such as RewardBench (Lambert et al., 2024) and RewardBench 2 (Malik et al., 2025) are the standard for testing preference generalization across diverse tasks. Their evaluation of safety is heavily focused on a model’s ability to refuse to generate harmful content. RewardBench, for example, draws from datasets like XSTest (Röttger et al., 2024) and Do-Not-Answer (Wang et al., 2023a), where the preferred response is a refusal and the rejected response is harmful. Similarly, RewardBench 2 leverages the CoCoNot taxonomy (Brahman et al., 2024) to assess compliance but explicitly excludes debatable categories to maintain a conservative stance on safety. Thus, while we recognize the importance of current benchmarks, their design may unintentionally steer models toward over-cautiousness, potentially sacrificing explanatory fidelity for broad-stroke safety in tasks that require nuanced discussion of sensitive topics.

To address this, **we introduce a benchmark focused on contextual accuracy in hate-explicit scenarios and propose a method to re-weight general-purpose, multi-attribute reward models to prioritize explanatory fidelity over blanket safety.**

3 HARM - A Hate-Aware Reward Model via Interpretable Attributes Re-weighting

3.1 Reward Model

Our RM is designed to score the explanatory fidelity of NLEs for hate speech. This model, which we call **Hate-Aware Reward Model (HARM)**, leverages pre-trained multi-attribute reward model by reweighting its interpretable dimensions.

Multi-Attribute Reward Model Backbone. Let x be a hate speech post from SBIC, and $E = G_i^{(t)}(x)$ be a model-generated explanation (NLE) produced by language model M_i under conditioning tier t , as defined in Section 4. Each explanation E is passed through a reward model that outputs a vector of d interpretable attribute scores: $\mathbf{v}_E \in \mathbb{R}^d$. These scores reflect semantic and stylistic properties such as helpfulness, coherence, truthfulness, and safety, depending on the training sources of the backbone reward model.

Mixture-of-Experts-Inspired Reward Modeling (HARM-MOE). We propose a RM formulation inspired by Mixture-of-Experts (MoE) architectures. Our approach conceptualizes the problem as a multitask learning challenge, similar to prior work in hate speech detection that uses specialized units to disentangle sentiment knowledge and improve system performance (Zhou et al., 2021). We hypothesize that an analogous division of labor is optimal for reward modeling in this domain, where the "tasks" correspond to judging explanations of either offensive or non-offensive content. To test this hypothesis, we designed an architecture that explicitly supervises two experts, one specializing in offensive content and the other in non-offensive content (*HARM-MOE-Off*). As a baseline, we compare it against a variant with the same two experts but without explicit supervision (*HARM-MOE-Un*), denoting the unsupervised setting.

Building on prior work (Christian et al., 2025; Jiang et al., 2025), we posit that the underlying reward function is inherently sparse: depending on the inputs nature, different evaluative criteria dominate. Leveraging this sparsity through conditional computation is increasingly seen as key to improving specialization, reducing interference, and scaling capacity efficiently (Pfeiffer et al., 2024; Shen et al., 2023; Du et al., 2022). We therefore adopt a fixed two-expert architecture:

Let $\mathbf{v}_E \in \mathbb{R}^d$ denote the interpretable attribute vector for an explanation E , where each element represents a different explanation feature (e.g., specificity, offensiveness, clarity). HARM computes a reward score using two expert branches: (i) a *positive expert (non-offensive)*, and (ii) a *negative expert (offensive)* that accounts for offensive aspects.

As highlighted in Figure 2, HARM-MOE-Off is designed in the following way:

(1) *Expert Gating.* Each expert uses a learned

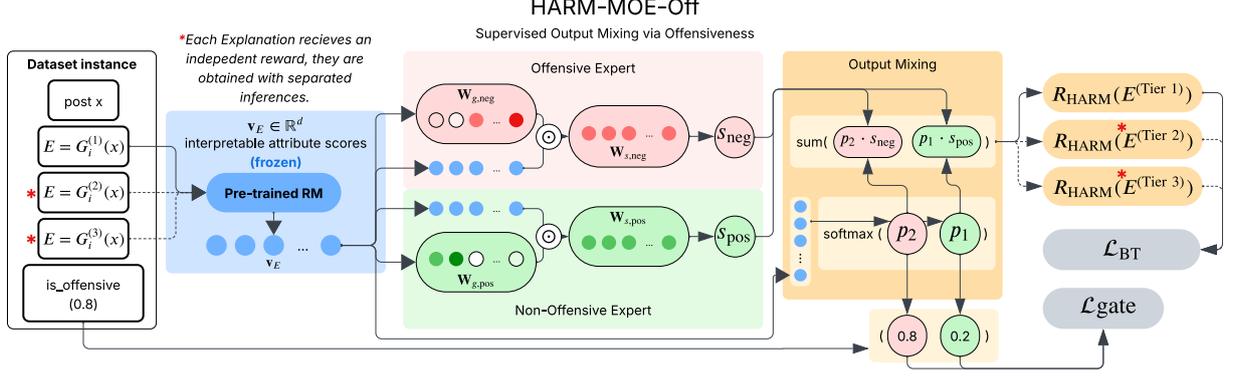


Figure 2: **HARM-MOE-Off**: Architecture inspired by Mixture of Experts with interpretable intermediate attribute re-weighting. The model receives one explanation candidate at a time and re-aggregates pre-trained attribute scores into a offensive-aware scalar reward.

gating matrix to select and scale relevant attributes:

$$\mathbf{g}_s = \sigma(\mathbf{W}_{g,s} \mathbf{v}_E) \odot \mathbf{v}_E, \quad s \in \{\text{pos}, \text{neg}\}$$

where $\mathbf{W}_{g,\text{pos}}, \mathbf{W}_{g,\text{neg}} \in \mathbb{R}^{d \times d}$ are learned gating matrices and $\sigma(\cdot)$ is the element-wise sigmoid. The result is a gated attribute vector for each expert.

(2) *Expert Scoring*. Each expert maps its gated vector to a scalar score using a learned projection:

$$s_{\text{pos}} = \mathbf{W}_{s,\text{pos}} \mathbf{g}_{\text{pos}}, \quad s_{\text{neg}} = \mathbf{W}_{s,\text{neg}} \mathbf{g}_{\text{neg}}$$

where $\mathbf{W}_{s,\text{pos}}, \mathbf{W}_{s,\text{neg}} \in \mathbb{R}^{1 \times d}$ are learned scoring weights.

(3) *Output Mixing*. A softmax over the input vector determines how much each expert contributes:

$$\mathbf{p} = \text{Softmax}(\mathbf{W}_p \mathbf{v}_E), \quad \text{with } \mathbf{p} = [p_1, p_2]$$

where $\mathbf{W}_p \in \mathbb{R}^{2 \times d}$ is a learned mixing matrix. The weights p_1 and p_2 represent the learned importance of each expert, conditioned on the explanations attributes.

(4) *Final Reward*. The final HARM score is a weighted combination of the expert scores:

$$R_{\text{HARM}}(E) = \beta \cdot (p_1 \cdot s_{\text{pos}} + p_2 \cdot s_{\text{neg}})$$

where β is a learned scalar used to calibrate reward scale.

We investigate two specialization mechanisms:

HARM-MOE-Un (Unsupervised Output Mixing): \mathbf{p} is trained using reward prediction alone, letting expert offensiveness specialization emerge naturally (Zoph et al., 2022).

HARM-MOE-Off (Supervised Output Mixing via Offensiveness): an auxiliary MSE loss $\mathcal{L}_{\text{gate}}$ aligns $\mathbf{W}_p \mathbf{v}_E$ with SBIC offensiveness labels, improving each expert offensiveness alignment (Zhou et al., 2022).

This approach integrates interpretable modular bias and task-aware supervision, relating to ensemble-inspired uncertainty modeling (Lakshminarayanan et al., 2017) and recent efforts in modular RL and evaluation (Aydn et al., 2025).

3.2 Training via BradleyTerry Loss

To optimize toward a correct multi-explanation score ranking, we adopt the BradleyTerry (Bradley and Terry, 1952) framework for pairwise preference learning, following recent advances in reward modeling (Rafailov et al., 2023; Wang et al., 2024a). Given two explanations (E_i, E_j), where E_i is preferred over E_j , the model learns to predict this preference through the scoring function R_{HARM} :

$$P(E_i \succ E_j | \theta) = \sigma(R_{\text{HARM}}(E_i) - R_{\text{HARM}}(E_j))$$

The main loss function is the BradleyTerry objective:

$$\mathcal{L}_{\text{BT}} = - \sum_{(i,j) \in \mathcal{P}} \log P(E_i \succ E_j | \theta)$$

where \mathcal{P} denotes the set of tiered explanation pairs.

Preference Pair Generation. We construct pairwise preferences by comparing explanation quality across tiers, conditional on the offensiveness of the post:

For posts marked as offensive (offensiveYN ≥ 0.5), stereotype information is crucial. Therefore, we generate full orderings: $E^{(\text{Tier } 3)} \succ E^{(\text{Tier } 2)} \succ E^{(\text{Tier } 1)}$

In these cases of *non-offensive post*, stereotype content is generally irrelevant or absent. The emphasis is on accurately conveying non-hatefulness. Consequently, Tier 2 and Tier 3 are equally preferred over Tier 1, as Tier 3 lacks human free-text annotations for non-offensive instances, and Tier 1 may inaccurately assess text offensiveness.: $E^{(\text{Tier } 3)} \sim E^{(\text{Tier } 2)} \succ E^{(\text{Tier } 1)}$

Importantly, we treat tier-based preferences as a *structural supervision signal* designed to probe a specific failure mode. Human validation shows the tendency holds on average but is noisy as show in Appendix C.

Total Training Objective. We define six loss components reflecting tiered comparisons between different content types (offensive vs. non-offensive), along with gating supervision for offensiveness:

$\mathcal{L}_{\text{off}}^{i>j}$: Bradley Terry (BT) loss comparing Tier j to Tier i for offensive content, where $(i, j) \in \{(1, 2), (1, 3), (2, 3)\}$.

$\mathcal{L}_{\text{non}}^{i>j}$: BT loss comparing Tier j to Tier i for non-offensive content, where $(i, j) \in \{(1, 2), (1, 3)\}$.

$\mathcal{L}_{\text{gate}}$: mean squared error (MSE) loss for output mixing supervision in the MOE-Off module.

The total loss is a weighted sum: $\mathcal{L}_{\text{total}} = \sum_k \lambda_k \mathcal{L}_k$, where each λ_k is a scalar (learned or manually tuned), and $\mathcal{L}_k \in \{\mathcal{L}_{\text{off}}^{1>2}, \mathcal{L}_{\text{off}}^{1>3}, \mathcal{L}_{\text{off}}^{2>3}, \mathcal{L}_{\text{non}}^{1>2}, \mathcal{L}_{\text{non}}^{1>3}, \mathcal{L}_{\text{gate}}\}$.

Additional appendices (G, H and I) test the impact of different strategies for losses combination.

4 SBIC-Explain - A Multi-Tiered synthetically generated NLE Dataset.

To create a dataset that allows for a controlled study of increasing hate-explicit information in NLE generation, the Social Bias Inference Corpus (SBIC) (Sap et al., 2020) was used as a foundation. For each post x in SBIC, we generate three tiers of explanations using a set of K pretrained language models, $\mathcal{M} = \{M_1, \dots, M_K\}$. The generation for model M_i is defined as: $G_i^{(t)}(x) = M_i(\text{Prompt}(x, \mathcal{C}^{(t)}))$, $t \in \{1, 2, 3\}$ where $\mathcal{C}^{(t)}$ is the conditioning set for each tier: $\mathcal{C}^{(1)} = \{\}$, $\mathcal{C}^{(2)} = \{\text{offensiveYN}\}$, $\mathcal{C}^{(3)} = \{\text{offensiveYN}, \text{targetMinority}, \text{targetCategory}, \text{targetStereotype}\}$. Prompts used to generate inference can be found in Appendix B.

To ensure that differences in explanations are due to the conditioning information and not sampling randomness, we use greedy decoding (i.e., "top_k=1") for all generations.

We process the official train/dev/test splits of SBIC, selecting only instances that include a stereotype annotation (i.e., with non-null values for targetMinority, targetCategory, and targetStereotype). To ensure a balanced representation of offensive and non-offensive content, we apply a downsampling strategy. Specifically, we retain all offensive instances ($\text{offensiveYN} \geq 0.5$) that contain semantic stereotype information and randomly downsample the non-offensive subset to match this distribution. This rebalancing yields the final **SBIC-Explain** dataset, comprising 30,899 aggregated instances.

To synthesize the dataset, models were selected taking as a reference the top four models of Open LLM Leaderboard² under 30B parameters and only official providers, using also their official GGUF quantized versions for efficient inference: *Qwen/Qwen3-14B*: A model by Alibaba Cloud (q_8 GGUF official quantization (Yang et al., 2025)³). *Microsoft/Phi-4*: A recent state-of-the art 14B parameter model from Microsoft (q_8 GGUF official quantization (Abdin et al., 2024)⁴). *TIIUAE/Falcon-3-10B-Instruct*: A high-performing instruction-tuned model from the Technology Innovation Institute (q_8 GGUF official quantization (Team, 2024)⁵). *Google/Gemma-3-27B-it*: A recent, powerful instruction-tuned model from Google (q4_0 GGUF official quantization (Team et al., 2025)⁶).

The final dataset was generated with the mentioned 30,899 unique SBIC posts, each post generated *three levels of NLE*, and each level of explanation was generated by *four models*, producing a dataset with a total of 370,788 synthetical NLEs. The final dataset version is grouped per post and per model, so that each instance has 3 NLEs that were produced by model M_i . The final dataset therefore has a total of 123,596 instances⁷. Additional experiments validating the tiered order with a human

²<https://huggingface.co/spaces/open-llm-leaderboard/open-llm-leaderboard/>

³<https://huggingface.co/Qwen/Qwen3-14B-GGUF>

⁴<https://huggingface.co/microsoft/phi-4-gguf>

⁵<https://huggingface.co/tiiuae/Falcon3-10B-Instruct-GGUF>

⁶https://huggingface.co/google/gemma-3-27b-it-qat-q4_0-gguf

⁷Available on supplementary material during revision

study on a sample dataset (Section 4.1), analyzing tiers and models diversity (Appendix E) and each tier text complexity (Appendix F) are provided.

4.1 SBIC-Explain Tier Human Validation.

We conducted a crowdsourced human evaluation on 200 instances (sampled with Latin Hypercube), each receiving 3 annotations, to test whether greater contextual guidance leads to higher *explanatory fidelity to the post* in generated explanations. Annotators on Amazon Mechanical Turk assessed post offensiveness and ranked the three tiered explanations for each instance according to explanatory fidelity (see Figure 3 for the template).

We collected rankings (1 to 3) for each tiered explanation, where lower mean ranks indicate higher preference (T1 = 2.18, T2 = 1.99, T3 = 1.82). A Friedman test revealed significant differences across tiers ($\chi^2 = 21.38$, $p = 2.27 \times 10^{-5}$), and one-sided Wilcoxon signed-rank tests confirmed that, in aggregate, higher-context tiers are preferred (Table 1).

Because tier preferences are not universal at the instance level, especially for subjective cases and for non-offensive posts where Tier 2 and Tier 3 can be similar, we additionally report *win-rates* (proportions) alongside significance tests: for each tier pair (e.g., $T3 > T2$), the win-rate quantifies the fraction of instances where annotators preferred one tier over the other, stratified by offensiveness. These proportions make the noise level in tier-based tendencies explicit and contextualize later automatic evaluations that use tier comparisons as a diagnostic probe rather than as error-free ground truth.

The complete experimental design and results are provided in Appendix C.

Direction ($T_i > T_j$)	Win-rate		Wilcoxon p-value
	Non-off. (< 0.5)	Off. (≥ 0.5)	All (one-sided)
2 > 1	0.64	0.48	0.008
3 > 1	0.63	0.60	$2.36e^{-6}$
3 > 2	0.52	0.59	0.014
3 > 2 > 1	0.26	0.21	–

Table 1: Human validation results on SBIC-Explain. We report (i) class-conditional win-rates (proportion of instances where annotators preferred one tier over another) and (ii) one-sided Wilcoxon signed-rank tests on tier rankings (lower rank = preferred).

Baseline RM	Tier	Accuracy
LaaJ*	2 > 1	0.19
	3 > 1	0.18
	3 > 2	0.11
URM (Lou et al., 2025)	2 > 1	0.58
	3 > 1	0.44
	3 > 2	0.23
ArmoRM (Wang et al., 2024a)	2 > 1	0.62
	3 > 1	0.49
	3 > 2	0.28
HARM MOE-Off	2 > 1	0.73
	3 > 1	0.87
	3 > 2	0.79

Table 2: SBIC-Explain pair-wise tiers accuracy of LaaJ, RMs and HARM for different tiers. *Mean across all models used as judges (Falcon, Gemma, Phi, Qwen), more detailed analysis on Appendix L)

5 Experiments and Results

To evaluate generalizability, we assess HARM on our custom SBIC-Explain dataset, which includes human-annotated rationales enabling the generation of Tier 1, 2, and 3 outputs, and the external Google Civil Comments dataset (Borkan et al., 2019), which supports only Tier 1 and 2 outputs due to its offensive-related labels (Appendix K). We used *ArmoRM-Llama3-8B-v0.1* (Wang et al., 2024a) as HARM’s backbone, prompt details are in Appendix B.4. All models were trained using the Adam optimizer with a learning rate of $2.5e-3$. To assess architectural trade-offs, we also test Linear and MLP variants as baselines (Appendix N).

5.1 Baseline Misalignment: Off-the-Shelf RMs Fail to Value Context.

Our results are presented based on two categories of methods: Reward models, utilizing ArmoRM-Llama3-8B-v0.1 (Wang et al., 2024a)⁸ and URM-LLaMa-3.1-8B (Lou et al., 2025)⁹, as-well as LaaJ based method. Specifically, the LaaJ based method used the same four models used for Synthetic NLE generation (Qwen, Phi, Falcon and Gemma, as detailed in Section 4) as judges. The rubric-based prompt templates was inspired by Prometheus (Kim et al., 2024) and can be further understood by

⁸<https://huggingface.co/RLHFFlow/ArmoRM-Llama3-8B-v0.1>

⁹<https://huggingface.co/LxzGordon/URM-LLaMa-3.1-8B>

Metric	Max Corr. (ArmoRM)	Max Corr. (URM)	Max Corr. (LaaJ)
Non-offensive	0.45 (Tier 3)	0.37 (Tier 3)	0.13 Tier 3)
Hate speech (not hate)	0.25 (Tier 2)	0.13 (Tier 3)	0.05 Tier 3)
Negative sentiment	-0.40 (Tier 3)	-0.43 (Tier 3)	-0.19 Tier 3)

Table 3: Correlations between reward and offensiveness related metric. More details in *Appendix D*.

reading *Appendix L*.

The results in *Table 2* confirm a critical failure: both state-of-the-art RMs and LaaJ based method prefer the same label-aware (Tier 2) explanations and penalize the inclusion of rich stereotype information (Tier 3). This demonstrates that general-purpose reward assignments are not suited for the nuances of this domain and actively discourage the generation of more insightful (but possibly more language sensitive) explanations.

Misalignment Hypothesis Rationale. We analyzed whether reward based on both LaaJ and RMs implicitly penalize sensitive or emotionally charged language by correlating their rewards with pretrained detectors for hate speech, offensiveness, and sentiment. The results from *Table 3* show that rewards are moderately *positively correlated with non-offensive content* and *negatively correlated with negative sentiment*, indicating that RMs tend to favor safer or sanitized language, even when explanations may require acknowledging offensive content.

Further analysis (*Appendix R*) shows that sanitized explanations received on average *37% higher rewards*, reinforcing that current RMs undervalue fidelity but socially grounded explanations.

5.2 HARM Performance: Correcting the Misalignment

We evaluate the explanations using our proposed *HARM-MoE-Off*. As shown in *Table 2*, HARM better matches the *tier-based preference probe* and the *human directional preference tendency* of explanation fidelity, assigning the highest rewards to Tier 3, demonstrating the relevance of re-weighting existing interpretable reward attributes. Notably, attributes like *coherence and complexity* receive higher weights, reflecting the importance of consistency and reasoning depth. A full analysis of learned weights is provided in *Appendix J*, with a dedicated discussion on text complexity in *Appendix F*, and practical examples in *Appendix R*.

Tier	Expert	MOE-Off	MOE-Un
1	Non-Offensive	-0.59	-0.48
1	Offensive	0.59	0.48
2	Non-Offensive	-0.78	-0.68
2	Offensive	0.78	0.68
3	Non-Offensive	-0.82	-0.73
3	Offensive	0.82	0.73

Table 4: Correlation of model **output mixing** ($\mathbf{W}_p \in \mathbb{R}^{2 \times d}$) and offensive ground truth, across model tiers.

Output Mixing: Offensive Post Prediction Accuracy. To evaluate whether supervised output mixing improves alignment between experts and offensive content, we analyzed the correlation between gating values and the SBIC `is_offensive` continuous annotations on the test set.

As shown in *Table 4*, the output mixing weights, particularly those selecting between the non-offensive and offensive experts, show strong and systematic correlations with the offensiveness signal. This alignment is further amplified in the HARM-MOE-Off setting, where gating is explicitly supervised using `is_offensive` annotations.

Despite receiving no explicit guidance, the unsupervised variant (HARM-MOE-Un) naturally develops gating behaviors aligned with offensive content, indicating that the emergence of expert specialization along offensive/non-offensive dimensions is an inductive bias reinforced by the task and data. Notably, the supervised model (HARM-MOE-Off) further amplifies this effect, achieving, on average, a 16.6% improvement in expert correlation with offensive content, demonstrating its ability to disentangle semantically meaningful behaviors more clearly and robustly. These findings highlight the interpretability of our output-gated architecture and point toward the broader potential of supervised and hybrid approaches for learning structured, disentangled representations of social biases and linguistic toxicity in a modular fashion.

6 Conclusion

We addressed a critical gap in the evaluation of Natural Language Explanations (NLEs) for hate speech by revealing how current reward models (RMs) as well as LLM-as-a-Judge (LaaJ) based method penalize contextually appropriate explanations, thereby misaligned with the goal of explanatory fidelity to the post in socially sensitive explanation.

To address this challenge, we introduce **SBIC-Explain**, a dataset comprising 370,788 LLM-generated explanations for offensive content, annotated across three tiers of increasing contextual richness: Tier 1 (text-only), Tier 2 (classification-aware), and Tier 3 (semantics-informed). This design is grounded in the hypothesis that explanations incorporating richer, stereotype-informed context yield higher perceived explanatory fidelity. We *validate this hypothesis through a human study* demonstrating that (i) the tiers are ranked significantly differently overall, and (ii) pairwise comparisons reveal significant differences across tiers, with Tier 1 perceived with lower explanatory fidelity than Tier 3.

We also presented **HARM**, a Hate-Aware Reward Model that reweights interpretable attributes to prioritize domain-relevant explanatory quality. HARM not only recovers the intended hierarchy across explanation tiers but also provides fine-grained interpretability, making it a practical tool for sensitive evaluation tasks. Our work aims to serve as a benchmark and inspiration for developing systems producing explanations with higher explanatory fidelity in sensitive contexts.

Limitations

Dependence on the backbone reward model. HARMs performance and learned weights depend on the choice of the backbone multi-attribute reward model (e.g., ArmoRM (Wang et al., 2024a)). Since such backbones are trained on broad, safety-oriented data, they may encode biases (e.g., toward refusals, aversion to sensitive language, or stylistic heuristics) that can propagate into HARMs scoring behavior, potentially limiting alignment in hate/offense explanation settings.

Synthetic explanations and limited human validation. Our evaluation relies on LLM-generated explanations conditioned on SBIC annotations rather than fully human-authored explanations. While this enables scale and controlled tiering, synthetic NLEs may omit the nuance, diversity, and implicit reasoning present in real moderation justifications. We partially mitigate this through a human validation study, but this study remains limited in size (200 instances) and exhibited substantial annotation noise: although 600 annotations were collected, many were discarded due to quality-control filters. Given the subjectivity of offensiveness and explanation evaluation, further studies with larger

and more diverse annotator pools are needed to consolidate the observed preference tendencies.

Scope and downstream use. HARM is currently designed as a post hoc evaluator and does not directly shape generation during training. Moreover, the extent to which the proposed methodology transfers to other high-stakes, linguistically sensitive domains (e.g., legal (Chalkidis et al., 2021) or medical (Agrawal et al., 2022; Han et al., 2025)) is untested and requires empirical validation. Integrating hate-aware attribute re-weighting into reinforcement learning or DPO-style training loops is a promising direction for future work, but it may introduce additional safety and robustness considerations.

SBIC annotation bias and potential amplification. Finally, our work inherits limitations from SBIC (Sap et al., 2020). Its offensiveness labels and stereotype fields reflect annotator subjectivity, and may embed systematic biases documented in prior analyses of toxicity and hate-speech annotation (Arhin et al., 2021; Conia et al., 2024). Because HARM is trained with tier-derived preference signals that partially depend on SBIC annotations, preference learning may *amplify* annotation artifacts: the model can learn to reward explanation patterns that align with SBICs annotator norms even when those norms do not generalize. While we reduce reliance on absolute labels by primarily learning within-instance tier comparisons and we also provide an unsupervised specialization variant, HARM should not be interpreted as a normative arbiter of what is offensive. SBIC-Explain and HARM are best viewed as tools for analyzing evaluator behavior under a specific annotation regime, and any deployment-facing use would require additional validation.

Acknowledgments

This work was partially supported by Capes and CNPq, project n. 441610/2023-4.

References

- Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa, Olli Saarikivi, and 8 others. 2024. [Phi-4 technical report](#). *Preprint*, arXiv:2412.08905.
- Monica Agrawal, Stefan Hegselmann, Hunter Lang, Yoon Kim, and David Sontag. 2022. [Large language models are few-shot clinical information extractors](#). *Preprint*, arXiv:2205.12689.
- Kofi Arhin, Ioana Baldini, Dennis Wei, Karthikeyan Natesan Ramamurthy, and Moninder Singh. 2021. [Ground-truth, whose truth? – examining the challenges with annotating toxic text datasets](#). *Preprint*, arXiv:2112.03529.
- Muhammad Umair Arshad and Waseem Shahzad. 2024. [Understanding hate speech: the HateInsights dataset and model interpretability](#). 10:e2372.
- Hüseyin Aydn, Kevin Godin-Dubois, Libio Goncalvez Braz, Floris den Hengst, Kim Baraka, Mustafa Mert Çelikok, Andreas Sauter, Shihan Wang, and Frans A. Oliehoek. 2025. [Sharpie: A modular framework for reinforcement learning and human-ai interaction experiments](#). *Preprint*, arXiv:2501.19245.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, and 12 others. 2022. [Training a helpful and harmless assistant with reinforcement learning from human feedback](#). *Preprint*, arXiv:2204.05862.
- Francesco Barbieri, Jose Camacho-Collados, Leonardo Neves, and Luis Espinosa-Anke. 2020. [Tweeteval: Unified benchmark and comparative evaluation for tweet classification](#). *Preprint*, arXiv:2010.12421.
- Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2019. [Nuanced metrics for measuring unintended bias with real data for text classification](#). *CoRR*, abs/1903.04561.
- Ralph Allan Bradley and Milton E. Terry. 1952. [Rank analysis of incomplete block designs: I. the method of paired comparisons](#). *Biometrika*, 39:324.
- Faeze Brahman, Sachin Kumar, Vidhisha Balachandran, Pradeep Dasigi, Valentina Pyatkin, Abhilasha Ravichander, Sarah Wiegrefe, Nouha Dziri, Khyathi Chandu, Jack Hessel, Yulia Tsvetkov, Noah A. Smith, Yejin Choi, and Hannaneh Hajishirzi. 2024. [The art of saying no: Contextual noncompliance in language models](#). *Preprint*, arXiv:2407.12043.
- Marc Brysbaert. 2019. [How many words do we read per minute? a review and meta-analysis of reading rate](#). *Journal of Memory and Language*, 109:104047.
- Ilias Chalkidis, Manos Fergadiotis, and Ion Androutsopoulos. 2021. [MultiEURLEX - a multi-lingual and multi-label legal document classification dataset for zero-shot cross-lingual transfer](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6974–6996, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Lichang Chen, Chen Zhu, Davit Soselia, Jiu-hai Chen, Tianyi Zhou, Tom Goldstein, Heng Huang, Mohammad Shoeybi, and Bryan Catanzaro. 2024a. [Odin: Disentangled reward mitigates hacking in rlhf](#). *Preprint*, arXiv:2402.07319.
- YanJun Chen, Dawei Zhu, Yirong Sun, Xinghao Chen, Wei Zhang, and Xiaoyu Shen. 2024b. [The accuracy paradox in RLHF: When better reward models don't yield better language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2980–2989, Miami, Florida, USA. Association for Computational Linguistics.
- Brian Christian, Hannah Rose Kirk, Jessica A.F. Thompson, Christopher Summerfield, and Tsvetomira Dumbalska. 2025. [Reward model interpretability via optimal and pessimal tokens](#). In *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency*, FAccT 25, page 10481059. ACM.
- Jaymari Chua, Yun Li, Shiyi Yang, Chen Wang, and Lina Yao. 2024. [Ai safety in generative ai large language models: A survey](#). *Preprint*, arXiv:2407.18369.
- Simone Conia, Daniel Lee, Min Li, Umar Farooq Minhas, Saloni Potdar, and Yunyao Li. 2024. [Towards cross-cultural machine translation with retrieval-augmented generation from multilingual knowledge graphs](#). *Preprint*, arXiv:2410.14057.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. [Ultrafeedback: Boosting language models with scaled ai feedback](#). *Preprint*, arXiv:2310.01377.
- Luigi Daniele and Suphavadeeprasit. 2023. [Amplify-instruct: Synthetically generated diverse multi-turn conversations for efficient llm training](#). *arXiv preprint arXiv:(coming soon)*.
- Zahra Delbari, Nafise Sadat Moosavi, and Mohammad Taher Pilehvar. 2024. [Spanning the spectrum of hatred detection: A persian multi-label hate speech dataset with annotator rationales](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17889–17897.

- Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, and 8 others. 2022. [Glam: Efficient scaling of language models with mixture-of-experts](#). *Preprint*, arXiv:2112.06905.
- Cesaire J. K. Fouodo, Lea L. Kronziel, Inke R. König, and Silke Szymczak. 2025. [Effect of hyperparameters on variable selection in random forests](#). *Preprint*, arXiv:2309.06943.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2024. [GPTScore: Evaluate as you desire](#). In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 6556–6576, Mexico City, Mexico. Association for Computational Linguistics.
- Tianyu Han, Lisa C. Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexei Figueroa, Alexander Löser, Daniel Truhn, and Keno K. Bressen. 2025. [Medalpaca – an open-source collection of medical conversational ai models and training data](#). *Preprint*, arXiv:2304.08247.
- Phu Gia Hoang, Canh Duc Luu, Khanh Quoc Tran, Kiet Van Nguyen, and Ngan Luu-Thuy Nguyen. 2023. [ViHOS: Hate speech spans detection for Vietnamese](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 652–669, Dubrovnik, Croatia. Association for Computational Linguistics.
- Sanne Hoeken, Sina Zarrieß, and Özge Alacam. 2024. [Hateful word in context classification](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 172–186, Miami, Florida, USA. Association for Computational Linguistics.
- Dirk Hovy, Taylor Berg-Kirkpatrick, Ashish Vaswani, and Eduard Hovy. 2013. [Learning whom to trust with MACE](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1120–1130, Atlanta, Georgia. Association for Computational Linguistics.
- Fan Huang, Haewoon Kwak, and Jisun An. 2023a. [Chain of explanation: New prompting method to generate quality natural language explanation for implicit hate speech](#). In *Companion Proceedings of the ACM Web Conference 2023*, WWW '23 Companion, page 9093, New York, NY, USA. Association for Computing Machinery.
- Fan Huang, Haewoon Kwak, Kunwoo Park, and Jisun An. 2024. [Chatgpt rates natural language explanation quality like humans: But on which scales?](#) *Preprint*, arXiv:2403.17368.
- Hui Huang, Xingyuan Bu, Hongli Zhou, Yingqi Qu, Jing Liu, Muyun Yang, Bing Xu, and Tiejun Zhao. 2025. [An empirical study of LLM-as-a-judge for LLM evaluation: Fine-tuned judge model is not a general substitute for GPT-4](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 5880–5895, Vienna, Austria. Association for Computational Linguistics.
- Olivia Huang, Eve Fleisig, and Dan Klein. 2023b. [Incorporating worker perspectives into MTurk annotation practices for NLP](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1010–1028, Singapore. Association for Computational Linguistics.
- Md Saroar Jahan and Mourad Oussalah. 2021. [A systematic review of hate speech automatic detection using natural language processing](#). *Preprint*, arXiv:2106.00742.
- Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan, Chi Zhang, Ce Bian, Chi Zhang, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. 2023. [Beavertails: Towards improved safety alignment of llm via a human-preference dataset](#). *Preprint*, arXiv:2307.04657.
- Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao Huang, Bill Yuchen Lin, and Wenhui Chen. 2024. [TIGER-Score: Towards building explainable metric for all text generation tasks](#). *Transactions on Machine Learning Research*.
- Junqi Jiang, Tom Bewley, Saumitra Mishra, Freddy Lecue, and Manuela Veloso. 2025. [Interpreting language reward models via contrastive explanations](#). In *The Thirteenth International Conference on Learning Representations*.
- Masahiro Kaneko and Danushka Bollegala. 2021. [Unmasking the mask – evaluating social biases in masked language models](#). *Preprint*, arXiv:2104.07496.
- Alex Kendall, Yarin Gal, and Roberto Cipolla. 2018. [Multi-task learning using uncertainty to weigh losses for scene geometry and semantics](#). *Preprint*, arXiv:1705.07115.
- Ahrii Kim. 2025. [RUBRIC-MQM : Span-level LLM-as-judge in machine translation for high-end models](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 6: Industry Track)*, pages 147–165, Vienna, Austria. Association for Computational Linguistics.
- Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024. [Prometheus: Inducing fine-grained evaluation capability in language models](#). In *The Twelfth International Conference on Learning Representations*.
- Youngwook Kim, Shinwoo Park, Youngsoo Namgoong, and Yo-Sub Han. 2023. [ConPrompt: Pre-training](#)

- a language model with machine-generated data for implicit hate speech detection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10964–10980, Singapore. Association for Computational Linguistics.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 64056416, Red Hook, NY, USA. Curran Associates Inc.
- Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2024. [Rewardbench: Evaluating reward models for language modeling](#). *Preprint*, arXiv:2403.13787.
- Yongqi Leng, Renren Jin, Yue Chen, Zhuowen Han, Ling Shi, Jianxiang Peng, Lei Yang, Juesi Xiao, and Deyi Xiong. 2025. [Praetor: A fine-grained generative LLM evaluator with instance-level customizable evaluation criteria](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10386–10418, Vienna, Austria. Association for Computational Linguistics.
- Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024. [Llms-as-judges: A comprehensive survey on llm-based evaluation methods](#). *Preprint*, arXiv:2412.05579.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. [G-eval: Nlg evaluation using gpt-4 with better human alignment](#). *Preprint*, arXiv:2303.16634.
- Xingzhou Lou, Dong Yan, Wei Shen, Yuzi Yan, Jian Xie, and Junge Zhang. 2025. [Uncertainty-aware reward model: Teaching reward models to know what is unknown](#). *Preprint*, arXiv:2410.00847.
- Lu Lu, Nathan Neale, Nathaniel D. Line, and Mark Bonn. 2022. [Improving data quality using amazon mechanical turk through platform setup](#). *Cornell Hospitality Quarterly*, 63(2):231–246.
- Qing Lyu, Marianna Apidianaki, and Chris Callison-Burch. 2024. [Towards faithful model explanation in NLP: A survey](#). *Computational Linguistics*, 50(2):657–723.
- Jitendra Singh Malik, Hezhe Qiao, Guansong Pang, and Anton van den Hengel. 2023. [Deep learning for hate speech detection: A comparative study](#). *Preprint*, arXiv:2202.09517.
- Saumya Malik, Valentina Pyatkin, Sander Land, Jacob Morrison, Noah A. Smith, Hannaneh Hajishirzi, and Nathan Lambert. 2025. [Rewardbench 2: Advancing reward model evaluation](#). *Preprint*, arXiv:2506.01937.
- Binny Mathew, Punyajoy Saha, Seid Muhie Yimam, Chris Biemann, Pawan Goyal, and Animesh Mukherjee. 2022. [Hatexplain: A benchmark dataset for explainable hate speech detection](#). *Preprint*, arXiv:2012.10289.
- Michael D McKay, Richard J Beckman, and William J Conover. 2000. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 42(1):55–61.
- Douglas Montgomery and Cahyono St. 2022. *Design and Analysis of Experiments, 9th Edition*.
- Rachna Narula and Poonam Chaudhary. 2024. [A comprehensive review on detection of hate speech for multi-lingual data](#). *Social Network Analysis and Mining*, 14(1):244.
- Huy Nghiem and Hal Daumé Iii. 2024. [HateCOT: An explanation-enhanced dataset for generalizable offensive speech detection via large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 5938–5956, Miami, Florida, USA. Association for Computational Linguistics.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). *Preprint*, arXiv:2203.02155.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, and 2 others. 2019. [Pytorch: An imperative style, high-performance deep learning library](#). *Preprint*, arXiv:1912.01703.
- John Pavlopoulos, Jeffrey Sorensen, Léo Laugier, and Ion Androutsopoulos. 2021. [SemEval-2021 task 5: Toxic spans detection](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*, pages 59–69, Online. Association for Computational Linguistics.

- Jonas Pfeiffer, Sebastian Ruder, Ivan Vuli, and Edoardo Maria Ponti. 2024. [Modular deep learning](#). *Preprint*, arXiv:2302.11529.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. [Direct preference optimization: Your language model is secretly a reward model](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 53728–53741. Curran Associates, Inc.
- Manikandan Ravikiran and Subbiah Annamalai. 2021. [DOSA: Dravidian code-mixed offensive span identification dataset](#). In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*, pages 10–17, Kyiv. Association for Computational Linguistics.
- Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. 2024. [Xstest: A test suite for identifying exaggerated safety behaviours in large language models](#). *Preprint*, arXiv:2308.01263.
- Ananya B. Sai, Akash Kumar Mohankumar, and Mitesh M. Khapra. 2020. [A survey of evaluation metrics used for nlg systems](#). *Preprint*, arXiv:2008.12009.
- Isadora Salles, Francielle Vargas, and Fabrício Benvenuto. 2025. [HateBRXplain: A benchmark dataset with human-annotated rationales for explainable hate speech detection in Brazilian Portuguese](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 6659–6669, Abu Dhabi, UAE. Association for Computational Linguistics.
- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. [The risk of racial bias in hate speech detection](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678, Florence, Italy. Association for Computational Linguistics.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. [Social bias frames: Reasoning about social and power implications of language](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5477–5490, Online. Association for Computational Linguistics.
- Chantal Shaib, Joe Barrow, Jiuding Sun, Alexa F. Siu, Byron C. Wallace, and Ani Nenkova. 2025. [Standardizing the measurement of text diversity: A tool and a comparative analysis of scores](#). *Preprint*, arXiv:2403.00553.
- Sheng Shen, Le Hou, Yanqi Zhou, Nan Du, Shayne Longpre, Jason Wei, Hyung Won Chung, Barret Zoph, William Fedus, Xinyun Chen, Tu Vu, Yuexin Wu, Wuyang Chen, Albert Webson, Yunxuan Li, Vincent Zhao, Hongkun Yu, Kurt Keutzer, Trevor Darrell, and Denny Zhou. 2023. [Mixture-of-experts meets instruction tuning: a winning combination for large language models](#). *Preprint*, arXiv:2305.14705.
- Tianxiang Sun, Junliang He, Xipeng Qiu, and Xuanjing Huang. 2022. [BERTScore is unfair: On social bias in language model-based metrics for text generation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3726–3739, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Falcon-LLM Team. 2024. [The falcon 3 family of open models](#).
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, and 197 others. 2025. [Gemma 3 technical report](#). *Preprint*, arXiv:2503.19786.
- Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. 2024a. [Interpretable preferences via multi-objective reward modeling and mixture-of-experts](#). In *EMNLP*.
- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Lingpeng Kong, Qi Liu, Tianyu Liu, and Zhifang Sui. 2024b. [Large language models are not fair evaluators](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9440–9450, Bangkok, Thailand. Association for Computational Linguistics.
- Yuxia Wang, Haonan Li, Xudong Han, Preslav Nakov, and Timothy Baldwin. 2023a. [Do-not-answer: A dataset for evaluating safeguards in llms](#). *Preprint*, arXiv:2308.13387.
- Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert, Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, and Oleksii Kuchaiev. 2023b. [Helpsteer: Multi-attribute helpfulness dataset for steerm](#). *Preprint*, arXiv:2311.09528.
- Martin Weysow, Aton Kamanda, Xin Zhou, and Houari Sahraoui. 2024. [Codeultrafeedback: An llm-as-a-judge dataset for aligning large language models to coding preferences](#). *Preprint*, arXiv:2403.09032.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- An Yang, Anpeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,

- Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. [Qwen3 technical report](#). *Preprint*, arXiv:2505.09388.
- Yue Yu, Zhengxing Chen, Aston Zhang, Liang Tan, Chenguang Zhu, Richard Yuanzhe Pang, Yundi Qian, Xuwei Wang, Suchin Gururangan, Chao Zhang, Melanie Kambadur, Dhruv Mahajan, and Rui Hou. 2025. [Self-generated critiques boost reward modeling for language models](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 11499–11514, Albuquerque, New Mexico. Association for Computational Linguistics.
- Lining Zhang, Simon Mille, Yufang Hou, Daniel Deutsch, Elizabeth Clark, Yixin Liu, Saad Mahamood, Sebastian Gehrmann, Miruna Clinciu, Khyathi Raghavi Chandu, and João Sedoc. 2023. [A needle in a haystack: An analysis of high-agreement workers on MTurk for summarization](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14944–14982, Toronto, Canada. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with bert](#). *Preprint*, arXiv:1904.09675.
- Xianbing Zhou, Yang Yong, Xiaochao Fan, Ge Ren, Yunfeng Song, Yufeng Diao, Liang Yang, and Hongfei Lin. 2021. [Hate speech detection based on sentiment knowledge sharing](#). In *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)*, pages 7158–7166, Online. Association for Computational Linguistics.
- Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Y. Zhao, Andrew Dai, Zhifeng Chen, Quoc Le, and James Laudon. 2022. [Mixture-of-experts with expert choice routing](#). In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA. Curran Associates Inc.
- Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2020. [Fine-tuning language models from human preferences](#). *Preprint*, arXiv:1909.08593.
- Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. 2022. [St-moe: Designing stable and transferable sparse expert models](#). *Preprint*, arXiv:2202.08906.

A Explainable Hate Speech Datasets

Table 5 shows that most explainable hate speech datasets center around highlight-based rationales, where annotators mark specific words or spans perceived as offensive. This design, seen in datasets like HateXplain, SemEval2021, and HateBRXplain, prioritizes the surface detection of hateful expressions. However, it often neglects the deeper reasoning or hateful rationale, the underlying social stereotypes, historical marginalization, or implied harm, behind why a phrase is offensive.

HateCOT introduces a promising step forward by generating synthetic natural language explanations (NLEs) that leverage broader context. However, its approach directly injects all available metadata (e.g., target labels, hate categories, spans) into the prompt, leaving open the question of how each type of information impacts explanation quality. This conflation makes it difficult to isolate which elements meaningfully improve explanation coherence or alignment with human values.

In contrast, our work (*SBIC-Explain*) proposes a more structured investigation. By focusing on the semantic contribution of stereotype-informed free-text rationales, specifically, the implied statements from SBIC, we aim to disentangle how this deeper contextual information affects model-generated explanations. Moreover, we go beyond human-written ground truth by evaluating how these enriched explanations are perceived by reward models, offering insights into both performance and alignment.

Dataset	Year	Language(s)	Source main	Domain	Explainability	Annotation	Size
HateBRXplain (Salles et al., 2025)	2025	Portuguese (Brazilian)	Instagram (comments)		Span rationales (text)	Expert annotators	7,000
HateInsights (Arshad and Shahzad, 2024)	2024	Urdu (Roman Arabic)	Twitter, Facebook		Span rationales (word/sentence)	Student annotators	11,782
HateCOT (Nghiem and Daumé Iii, 2024)	2024	English	Multi \approx source (8 corpora)		Synthetic NLEs	GPT3.5-Turbo + human curation	52,137
PHate (Delbari et al., 2024)	2024	Persian (Farsi)	Twitter		Span rationales + target labels	Expert annotators	7,000
HateWiC (Hoecken et al., 2024)	2024	English	English Wiktionary		Span Rationales (text)	Three annotators	4,029
ViHOS (Hoang et al., 2023)	2023	Vietnamese	Facebook, YouTube		Span rationales (text)	Human annotators	11,056
HateXplain (Mathew et al., 2022)	2020	English	Twitter, Gab		Span rationales (text)	AMT crowdworkers	20,148
SemEval-2021 Task 5 (Toxic Spans) (Pavlopoulos et al., 2021)	2021	English	Civil Comments (Wikipedia talk)		Span rationales (text)	Crowd workers	10,629
DOSA (Ravikiran and Annamalai, 2021)	2021	Tamil English, Kannada English	YouTube comments		Span rationales (words)	Human annotators	4,786 (Tamil) / 1,097 (Kannada)
Social Bias Frames (SBIC) (Sap et al., 2020)	2020	English	Twitter, Reddit, Gab, extremist forums		Categorical labels + freetext implied statements	Crowdsourced (MTurk)	44,671 posts
SBIC-Explain (OURS)	2025	English	Twitter, Reddit, Gab, extremist forums		3-level Synthetic NLEs based on 4 different generative models	Crowdsourced (MTurk)	123,599 instances and 370,788 NLEs

Table 5: Overview of explainable hate speech detection datasets.

B Prompt Templates

Here we detail the exact prompt templates used for generating the NLEs for each tier. Different models had slight prompt variations to accommodate for different expected prompt structures (e.g. models that expect explicit reasoning tokens or slight variation in system/user/assistant chat tokens to specify who is talking), always following each models best practices.

B.1 Tier 1 Prompt

```

1 <|im_start|>user
2 ===== INFO TO FOLLOW WHEN
   GENERATING ANSWER
3 Post Text: {post}
4 ===== INFO TO FOLLOW WHEN
   GENERATING ANSWER<|im_end|>
5
6 <|im_start|>user
7 Generate a natural language
   explanation of why this is or is

```

```

8 not to be considerate offensive.
9 Generate the most concise explanation
   version possible while still
   explaining all important aspects.
10 You dont need to rewrite the whole
    sentence to explain it, but you can
    rewrite or reference specif parts
    to help on your explanation.<|
    im_end|>
11 Write your explanation in english.
12 <|im_end|>
13 <|im_start|>assistant /no_think
14
15 <think>
16
17 </think>
18
19 Explanation:

```

B.2 Tier 2 Prompt

```

1 <|im_start|>user
2 ===== INFO TO FOLLOW WHEN
   GENERATING ANSWER

```

```

3 Post Text: {post}
4 {'POST TEXT IS CONSIDERED OFFENSIVE '
   if is_offensive>=0.5 else 'POST
   TEXT IS NOT CONSIDERED OFFENSIVE '}
5 ===== INFO TO FOLLOW WHEN
   GENERATING ANSWER<|im_end|>
6
7 <|im_start|>user
8 Generate a natural language
   explanation of why this is or is
   not to be considerate offensive.
9 Generate the most concise explanation
   version possible while still
   explaining all important aspects.
10 You dont need to rewrite the whole
   sentence to explain it, but you can
   rewrite or reference specif parts
   to help on your explanation.<|
   im_end|>
11 Write your explanation in english.
12 <|im_end|>
13
14 <|im_start|>assistant /no_think
15
16 <think>
17
18 </think>
19
20 Explanation:

```

B.3 Tier 3 Prompt

```

1 <|im_start|>user
2 ===== INFO TO FOLLOW WHEN
   GENERATING ANSWER
3 Post Text: {post}
4 {'POST TEXT IS CONSIDERED OFFENSIVE '
   if is_offensive>=0.5 else 'POST
   TEXT IS NOT CONSIDERED OFFENSIVE '}
5 Targeted Minority: {targetMinority}
6 Targeted Category: {targetCategory}
7 Targeted Stereotype: {targetStereotype
   }
8 ===== INFO TO FOLLOW WHEN
   GENERATING ANSWER<|im_end|>
9
10 <|im_start|>user
11 Generate a natural language
   explanation of why this isor is
   not to be considerate offensive.
12 Generate the most concise explanation
   version possible while still
   explaining all important aspects.
13 You dont need to rewrite the whole
   sentence to explain it, but you can
   rewrite or reference specif parts
   to help on your explanation.<|
   im_end|>
14 Write your explanation in english.
15 <|im_end|>
16
17 <|im_start|>assistant /no_think
18
19 <think>
20
21 </think>
22
23 Explanation:

```

B.4 Templates For ArmoRM and URM

Here we detail the exact template used to generate Rewards using ArmoRM and URM. We followed the suggested structure as described in their Hugging Face pages¹⁰.

```

1 [{
2   "role": "user", "content": Generate a
   natural language explanation of why
   this post text is or is not
   considered offensive.
3 Post text: {<post_text>}
4 },
5 {"role": "assistant", "content": {<
   ai_generate_NLE>}
6 }]

```

¹⁰ArmoRM - <https://huggingface.co/RLHFlow/ArmoRM-Llama3-8B-v0.1>
 URM - <https://huggingface.co/LxzGordon/URM-LLaMa-3.1-8B>

C Detailed Human Evaluation

To test our hypothesis that richer contextual guidance yields higher perceived explanatory fidelity (Tier 3 > Tier 2 > Tier 1), we conducted a human evaluation on 200 representative instances subsampled from the **SBIC-Explain** dataset’s test split.

Following Design of Experiments (DOE) best practices (Montgomery and St, 2022), we avoided biases from random sampling (e.g., clustering in low-variance regions or over-representing edge cases) by using stratified, space-filling methods. This ensured even coverage across high-dimensional features, reducing confounders like model-specific or stereotype-skewed explanations and boosting generalizability.

We applied Latin Hypercube Sampling (LHS) (McKay et al., 2000), which divides the feature space into equal-probability strata and samples one point per stratum for maximal dispersion and low correlation, ideal for our nonlinear, high-dimensional data over simple random methods. For each test instance, we built a 3,074-dimensional vector with the: (i) normalized offensiveness score; (ii) LLM identifier (an integer for each unique LLM); and (iii) sentences embeddings based on all-mpnet-base-v2¹¹ for the post text and three generated explanations (2 + 4 × 768 dims). LHS generated 200 points in [0,1]³⁰⁷⁴. To align these with actual dataset instances, we firstly scaled the LHS generated features to each original feature original scale and used **k-NN** to retrieve the nearest test split instance per point, yielding a diverse subsample spanning offensiveness, models and texts.

The resulting subsample exhibited a balanced offensiveness score distribution (mean: 0.56, std: 0.39, range: [0,1]) and varied model representation (Qwen3-14B-Q8_0: 36.7%, Falcon3-10B-Instruct-Q8_0: 28.9%, Gemma-3-27b-it-Q4_0: 18.4%, Phi4-14B-Q8_0: 16.0%).

Annotations on this sample were crowdsourced via MTurk from Master-badge workers in English-dominant regions (USA, Canada, Ireland, UK, Australia) given the importance of have native english speaker to judge implicit. Each of 200 samples was evaluated by 3 annotators, that were fairly paid \$ 0.15 per task (\$7.25/hour U.S. minimum, based on average 1.5 min/task).

As shown in Figure 3, the interface displayed a post with three anonymized explanations (ran-

domly ordered as "A/B/C"). Tasks included:

Offensiveness Classification. Rate whether the original post is offensive using a 3-point ordinal scale (Yes this could be offensive/Maybe Im not sure/No this seems harmless)

Explanatory Fidelity. For each explanation, rate how faithfully it explains the post content using a 3-point ordinal scale (Accurate Clearly faithful/Partially accurate Somewhat faithful/Not accurate Not faithful)

Preference Ranking. Assign each explanation a unique rank (1-3, where 1 = best/most faithful) based on overall faithfulness. We defined faithfulness in the experiment as follows: *A faithful explanation is one that accurately points out the offensive content (for example, a stereotype, slur, or negative reference to a group) and clearly explains why it is offensive. Less faithful explanations may be vague, overlook the main issue, or introduce unsupported details.*

C.1 Quality Control and Filtering

Due to the prevalence of low-quality annotations commonly observed in crowdsourcing platforms (Huang et al., 2023b), we implemented a multi-stage quality control process:

Ranking Consistency: We required annotators to assign unique ranks to each explanation (no tied rankings). Annotations with repeated ranks were automatically flagged as invalid.

Content Engagement: We identified several patterns indicative of disengaged annotation:

- Ranking patterns (always ranking explanations as 1, 2, 3 or 3, 2, 1 regardless of content).
- Extremely short completion times inconsistent with thorough evaluation, calculated based on the average silent human reading speed (238 words per minute (wpm) (Brysbaert, 2019)) and the amount of words that had to be read.
- Intermediary explanation reasoning doesn’t match the final ranking. (e.g Explanation A was the only explanation marked as No faithful and yet occupies ranking position 1).

Annotators matching two or more of these criteria were disqualified.

¹¹<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

View full instructions

Post

#{postText}

Does this post address stereotypes or topics that could be offensive to any group of people? -- Select --

Explanation A:

#{explanationA}

Does this explanation accurately identify whether the post targets any group(s) or stereotype(s), and clearly explain why it may or may not be offensive?

-- Select --

Explanation B:

#{explanationB}

Does this explanation accurately identify whether the post targets any group(s) or stereotype(s), and clearly explain why it may or may not be offensive?

-- Select --

Explanation C:

#{explanationC}

Does this explanation accurately identify whether the post targets any group(s) or stereotype(s), and clearly explain why it may or may not be offensive?

-- Select --

Your Ranking

Assign each rank (1, 2, 3) to exactly one explanation.

Explanation A

Select rank ▼

Explanation B

Select rank ▼

Explanation C

Select rank ▼

Submit

Figure 3: Amazon Mechanical Turk (MTurk) HTML template for collecting human annotation.

Filtering Results. Our filtering process reduced the annotation set from 600 to 343 higher-quality annotations (57.1% retention rate). This filtering rate aligns with previous work on crowd sourced annotation quality (Hovy et al., 2013; Zhang et al., 2023), where substantial portions of collected annotations often require removal (discard rates up to 70%) due to quality concerns, and 20% to 30% invalid repos discards based on control mechanisms such as attention checks and constrained task completion time (Lu et al., 2022).

C.2 Results

Dataset Characteristics. Our final analysis considered 343 annotations, with an average of 1.72 annotations. The retained annotations were contributed by 26 annotators that analyzed 200 instances.

Hypothesis Validation. Our central hypothesis that explanations with higher contextual guidance would receive better rankings (Tier 3 > Tier 2 > Tier 1) was supported by the data. The mean rank-

ings for each tier (T1 = 2.18, T2 = 1.99 and T3 = 1.82), where lower values indicate better rankings, confirming the predicted preference order. A Friedman test, chosen for its suitability in comparing multiple related samples with ordinal ranking data (and non-parametric properties), revealed statistically significant differences across the three explanations tiers ($\chi^2 = 21.38, p = 2.27e^{-5}$), supporting our hypothesis that contextual guidance generates significantly different perceived explanatory fidelity.

C.2.1 Content-Type Analysis

We conducted separate pairwise analyses for offensive and non-offensive posts to examine whether content sensitivity affects explanation preferences. Table 6 reports one-sided Wilcoxon signed-rank tests, which are appropriate for paired ordinal rankings and do not assume normality. The one-sided direction tests whether the higher-context tier is preferred (e.g., whether Tier 3 is ranked better than Tier 2), matching our directional hypothesis.

In addition to significance tests, we also report

Content Type	Comparison	p-value	Win-rate
Offensive (≥ 0.5)	$T2 > T1$	0.56	0.48
	$T3 > T1$	$6.82e^{-3}$	0.60
	$T3 > T2$	$2.61e^{-3}$	0.59
	$T3 > T2 > T1$	–	0.21
Non-Offensive (< 0.5)	$T2 > T1$	$4.00e^{-4}$	0.64
	$T3 > T1$	$3.59e^{-5}$	0.63
	$T3 > T2$	0.39	0.52
	$T3 > T2 > T1$	–	0.26
All	$T2 > T1$	0.008	0.56
	$T3 > T1$	$2.36e^{-6}$	0.61
	$T3 > T2$	0.014	0.55
	$T3 > T2 > T1$	–	0.23

Table 6: One-sided Wilcoxon signed-rank tests on tier rankings (lower rank = preferred) and corresponding human win-rates (proportions) for pairwise preferences, stratified by offensiveness. $T3 > T2 > T1$ reports the proportion of instances where annotators preferred the strict tier ordering.

win-rates (proportions) for each comparison. For a pair such as $T3 > T2$, the win-rate is the fraction of evaluated instances where annotators preferred Tier 3 over Tier 2. These proportions are informative because tier preferences are noisy at the instance level (especially in subjective cases and for non-offensive posts, where Tier 2 and Tier 3 can be similar), and they make the strength of the directional tendency explicit beyond p-values.

Overall, the results indicate that richer contextual guidance (Tier 3) is preferred over minimal context (Tier 1) for both offensive and non-offensive posts. Consistent with the intuition that stereotype annotations are most informative for offensive instances, the Tier 3 versus Tier 2 difference is statistically significant only in the offensive subset, while it is weaker (and often indistinguishable) for non-offensive posts where Tier 3 lacks additional free-text stereotype content. The win-rates further contextualize this pattern by showing that these preferences are directional tendencies rather than universal instance-level orderings.

D Correlation of offensive related metrics and Reward

To better understand how RMs may implicitly disincentivize sensitive or emotionally charged content, we compute Pearson correlation coefficients between model-assigned rewards and external language classifiers targeting hate speech, offensiveness, and sentiment.

We use the following pretrained detectors from `cardiffnlp`: *Sentiment (sent)*: `twitter-roberta-base-sentiment` (Barbieri et al., 2020)¹². *Hate speech (hate)*: `twitter-roberta-base-hate`¹³. *Offensiveness (off)*: `twitter-roberta-base-offensive`¹⁴.

Overall, Table 7 shows moderate positive correlations between rewards and *non-offensive* content (up to $r = 0.45$ in Tier 3 for ArmoRM), suggesting that RMs may indeed favor less offensive or emotionally negative outputs. This aligns with previous findings that reward models often exhibit an implicit bias toward safe or sanitized language (Christian et al., 2025; Jiang et al., 2025).

Moreover, reward tends to be negatively correlated with negative sentiment scores, especially in more complex explanations, reaching values as low as $r = -0.43$ in Tier 3 (URM). This implies that emotionally charged or confrontational explanations, while potentially necessary in sensitive domains like hate speech, may be under-rewarded. These trends highlight the presence of systematic bias in reward modeling, further motivating the development of HARMs multi-expert mechanism.

In contrast, LaaJ exhibits weaker correlations across all tiers, with values ranging from -0.22 to 0.27, suggesting it is less influenced by content offensiveness or sentiment. For higher tiers (2 and 3), LaaJs correlations diminish further, highlighting its potential against the bias seen in other models.

To further investigate, we manually sanitized offensive explanations, removing references to stereotypes and bias, and compared their reward scores using ArmoRM (Appendix R). On average, sanitized explanations received *37% higher rewards*, supporting the hypothesis that current RMs under-value explanations with higher explanatory fidelity and socially grounded content.

¹²<https://huggingface.co/cardiffnlp/twitter-roberta-base-sentiment>

¹³<https://huggingface.co/cardiffnlp/twitter-roberta-base-hate>

¹⁴<https://huggingface.co/cardiffnlp/twitter-roberta-base-offensive>

Text	Reward Metric	ArmoRM (Wang et al., 2024a)			URM (Lou et al., 2025)			LaaJ*		
		Tier 1	Tier 2	Tier 3	Tier 1	Tier 2	Tier 3	Tier 1	Tier 2	Tier 3
Post	not hate	0.20	0.25	0.24	-0.07	-0.11	0.07	0.04	-0.07	-0.04
	not off	0.36	0.41	0.39	-0.01	0.18	0.20	-0.16	0.09	0.06
	sent (neg)	-0.20	-0.25	-0.24	-0.01	-0.14	-0.17	0.13	-0.07	-0.06
Tier 1	not hate	0.13	0.18	0.17	-0.07	0.07	0.08	-0.11	0.06	0.04
	not off	0.36	0.44	0.42	-0.05	0.21	0.25	-0.21	0.11	0.08
	sent (neg)	-0.31	-0.37	-0.38	-0.10	-0.27	-0.34	0.23	-0.09	-0.12
Tier 2	not hate	0.13	0.16	0.16	-0.05	0.05	0.10	-0.09	0.04	0.05
	not off	0.37	0.43	0.43	-0.03	0.17	0.27	-0.18	0.07	0.08
	sent (neg)	-0.25	-0.36	-0.40	0.05	-0.24	-0.40	0.27	-0.07	-0.16
Tier 3	not hate	0.13	0.17	0.17	-0.05	0.06	0.13	-0.10	0.04	0.05
	not off	0.34	0.42	0.45	-0.03	0.20	0.37	-0.22	0.06	0.13
	sent (neg)	-0.24	-0.35	-0.40	0.05	-0.25	-0.43	0.27	-0.06	-0.19

Table 7: Correlation between different methods reward results and sentiment, offensive and hate measurements of related texts (post, NLE 1,2,3 = Tier 1,2,3). *LaaJ: The mean reward off all models utilized as judges (Falcon, Gemma, Phi, Qwen) was utilized to calculate correlation.

E SBIC-Explain Diversity

To highlight the importance of considering different models, diversity of model-generated content was assessed using *compression ratio* of the concatenated outputs (Shaib et al., 2025). This metric, defined as the ratio of the compressed size (via gzip) to the original uncompressed size, serves as a proxy for output redundancy. Lower compression ratios indicate higher diversity, as the content contains fewer repeated substrings and patterns. We report each model’s sentence length (Avg Words) right-trimming sentences to a fixed size of 64 words (considering each word is separated by white space), because Qwen3 was the model that generated the smaller sentences, with an average of 64 words per sentence, and, as stated by (Shaib et al., 2025), the length of the analyzed text has to be reported alongside all these scores.

Table 8 reveals notable variations in sentence length across models and tiers of generated NLEs. Falcon produces the longest sentences, with an average of 142 words across tiers, peaking at 161 words in Tier 1. In contrast, Qwen generates the shortest sentences, averaging 64 words across tiers, with values (68, 62, and 61 words for Tiers 1, 2, and 3, respectively). Gemma exhibits variability, with an average of 133 words, showing a slight increase from 123 words in Tier 2 to 144 words in

Tier Model	Avg Words			Avg
	1	2	3	
Gemma	144	123	<u>136</u>	133
Falcon	<u>161</u>	<u>137</u>	124	142
Qwen	68	62	61	64
Phi	103	90	79	91

Table 8: Average Sentence length (words) of Generated NLEs of SBIC-Explain by Generative model and Tier. Unlines represent highest in row, while **bold** represent highest in column.

Tier 1. Phi maintains moderate sentence lengths, averaging 91 words, with a range from 79 words in Tier 3 to 103 words in Tier 1, reflecting a balanced output structure. These differences suggest that sentence length influences diversity, with shorter sentences potentially reducing repetitive patterns and longer sentences increasing the likelihood of structural or lexical redundancy.

From Table 9 results, Gemma-3 generates the most *diverse*, while Phi-4 and Qwen-3 tend to repeat structures or lexical patterns across samples. Falcon-3 occupies a middle ground, generating content with moderate variability.

Tiers Model	1	2	3	Avg	Avg Words
Gemma3	<i>3.17</i>	3.23	3.34	3.24	63
Falcon3	3.31	3.5	3.85	3.56	63
Qwen3	3.65	3.94	4.3	3.96	59
Phi4	3.66	3.92	4.49	4.02	63

Table 9: Compression ratios (as a proxy of diversity) of generated NLE tiers. *Italic* text highlight lowest in columns (more diverse ↓) and **bold** highlight highest in columns.

F Text Complexity Metrics

text Metric	Post	Tier 1	Tier 2	Tier 3
flesch_reading_ease ↓	72.998	29.349	29.154	25.107
mcalpine_eflaw ↑	20.526	29.12	28.495	28.402
gunning_fog ↑	8.689	17.782	17.668	18.37
smog_index ↑	8.196	15.924	15.797	16.276
automated_readability_index ↑	7.31	15.323	15.204	15.673
linsear_write_formula ↑	7.748	14.824	14.813	15.327
flesch_kincaid_grade ↑	6.667	14.392	14.329	14.953
coleman_liiau_index ↑	7.114	14.089	14.139	14.609
dale_chall_readability_score ↑	8.956	11.27	11.3	11.372
spache_readability ↑	4.46	7.035	6.997	7.115

Table 10: Complexity metrics (readability, complexity, and grade level). ↑ and ↓ indicate the relationship of the given metric and increase in complexity.

To further investigate the high importance of complexity-related attributes observed in HARMs learned weights, particularly *helpsteer-complexity* and *code-complexity*, we analyze textual complexity across different input and explanation types in our dataset.

We apply a suite of standard readability and complexity metrics using the `textstat` Python package¹⁵. These include traditional readability scores (e.g., Flesch Reading Ease), grade-level indicators (e.g., Gunning Fog Index, SMOG), and composite formulas (e.g., McAlpine EFLAW). Results are presented in Table 10, comparing complexity across the original post texts and the three tiers of generated explanations.

We observe a clear trend on Table 10: explanations, especially those from Tier 3, consistently exhibit higher complexity than the original posts or lower-tier explanations. For example, Tier 3 explanations have the lowest Flesch Reading Ease (25.1, lower is more complex) and the highest Gunning Fog Index (18.4, higher is more complex), indicating more sophisticated and cognitively demanding language. This aligns with our hypothesis that

more contextually informed explanations require greater linguistic and conceptual complexity.

These findings support the observed model behavior: attributes related to textual complexity receive higher weighting during HARMs scoring process, likely because they capture deeper reasoning and nuance, critical for high-quality hate speech explanation.

¹⁵<https://pypi.org/project/textstat/>

G MOE-Off - Losses Combinations (Losses Included or Excluded)

In this section, we investigate the effect of selectively including or excluding different training tiers when optimizing the MOE-Off reward model. Our goal is to understand how various data configurations contribute to generalization and performance across models, and which combinations offer the most effective supervision signal for the intended reward modeling objectives.

Table 11 presents pairwise accuracies, general accuracy, and reward score distance across multiple tier combinations and backbone models. We compare these configurations against our reference models (ArmoRM and URM), and also explore the impact of using all available data versus subsets grouped by tier.

Notably, training on Tiers 1 & 3 consistently yields the strongest results, achieving the highest overall accuracy (**0.807**) and the best average score distance (**0.111**). This indicates that the contrast between these two extremes (Tier 1 (low preference) and Tier 3 (high preference)) provides the clearest and most informative supervision signal. In contrast, combinations involving only Tiers 2 & 3 or Tiers 1 & 2 exhibit more moderate performance, likely due to the reduced contrast in reward preference, which makes learning signal less distinct.

Among backbone models, the **All** ensemble consistently ranks highest or very close to best, highlighting the benefit of incorporating diverse model perspectives during training. This supports our hypothesis that model diversity helps generalize the reward function across different architectures.

Interestingly, although Gemma and Qwen achieve strong results in some configurations, they also show higher variance, particularly Qwen, which exhibits significant volatility in score distance, likely due to instability in reward scaling. These findings suggest that while individual backbones can be powerful, their behavior must be carefully calibrated when used in isolation.

These results also align with diversity analysis (Appendix E) trends, where *HARM trained on data only generated by Falcon, Qwen, Phi, and Gemma* achieve progressively lower overall test accuracy, respectively, and training on all models synthetic data generated higher testing accuracy.

Interestingly, although Gemma-3 explanations are the most diverse, models trained alone on its data generalize the worst. This may suggest that high

diversity alone does not ensure effective reward model training. In contrast, Falcon-3’s intermediate diversity may offer a more balanced training signal, enabling stronger generalization across unseen model outputs.

Overall, this analysis reinforces two key insights: (1) contrastive supervision from clearly distinguishable preference tiers (especially Tiers 1 and 3) is critical for effective reward modeling, and (2) incorporating multiple models helps stabilize and improve reward quality across the board.

Model Configuration		Pairwise Accuracy				Reward Score Distance ($diff(i, j) = (i - j)/j$)			
Train Tiers	Train Model	2>1	3>1	3>2	Mean	2-1	3-1	3-2	Mean
Ref. Results	ArmoRM	0.616	0.489	0.282	0.524	0.039	0.005	-0.053	0.016
	URM	0.585	0.445	0.229	0.49	0.091	0.011	-0.136	0.038
Tiers 1 & 2	Falcon	0.747	0.67	0.385	0.669	0.335	0.318	-0.008	0.285
	Gemma	0.725	0.687	0.42	0.676	2.761	3.172	-0.283	2.883
	Phi	0.719	0.697	0.436	0.674	0.094	0.106	-0.008	0.09
	Qwen	0.743	0.755	0.537	0.718	0.131	0.164	0.027	0.137
	All	0.779	0.745	0.454	0.721	0.043	0.045	0.001	0.04
Tiers 1 & 3	Falcon	0.717	0.816	0.747	0.768	0.036	0.059	0.035	0.048
	Gemma	0.669	0.735	0.658	0.709	0.027	0.037	0.017	0.033
	Phi	0.644	0.816	0.8	0.743	0.03	0.098	0.081	0.067
	Qwen	0.655	0.821	0.815	0.756	0.019	0.052	0.042	0.038
	All	0.729	0.868	0.792	0.807	0.073	0.144	0.092	0.111
Tiers 2 & 3	Falcon	0.572	0.787	0.827	0.697	0.007	0.044	0.049	0.028
	Gemma	0.345	0.492	0.75	0.448	-0.044	0.006	0.084	-0.01
	Phi	0.479	0.738	0.847	0.632	-0.019	0.084	0.136	0.044
	Qwen	0.551	0.767	0.84	0.686	0.003	0.04	0.047	0.025
	All	0.495	0.763	0.86	0.655	-0.039	0.152	0.203	0.07
All Tiers	Falcon	0.687	0.813	0.76	0.75	0.068	0.196	0.2	0.142
	Gemma	0.575	0.707	0.744	0.654	0.006	0.02	0.026	0.014
	Phi	0.625	0.807	0.801	0.724	0.019	0.068	0.05	0.044
	Qwen	0.639	0.804	0.816	0.736	0.01	0.042	0.031	0.023
	All	0.689	0.849	0.808	0.775	0.03	0.083	0.069	0.059

Table 11: Performance of **HARM-MOE-Off** reward model across different training configurations. All reported results are based on the entire test set, only training is being altered.

H MOE-Off - Losses Combination (Grid Search 0, 0.5, 1.)

Loss Component	Importance	Correlation
$\mathcal{L}_{\text{non}}^{1>3}$	0.685	0.761
$\mathcal{L}_{\text{off}}^{1>2}$	0.106	0.129
$\mathcal{L}_{\text{off}}^{1>3}$	0.076	0.186
$\mathcal{L}_{\text{off}}^{2>3}$	0.052	0.080
$\mathcal{L}_{\text{non}}^{1>2}$	0.044	-0.128
$\mathcal{L}_{\text{gate}}$	0.037	0.019

Table 12: Feature importance and correlation analysis for MOE-Off loss weight parameters

To optimize the loss weighting strategy for our MOE-Off architecture, we conduct a comprehensive grid search across all loss components, system-

atically varying each weight parameter between 0, 0.5, and 1.0. This exploration allows us to identify the optimal combination of loss terms and understand the relative importance of different training objectives in our hate-aware reward modeling framework.

We employ a feature importance analysis technique inspired by (Fouodo et al., 2025), where we train a random forest model using hyperparameter configurations as inputs and place **accuracy** as target outputs. The random forest’s feature importance values reveal which loss weights contribute most significantly to model performance, while correlation analysis shows the direction and strength of these relationships.

Table 12 presents the feature importance and correlation analysis for each loss weight paramete-

Test Acc.	$\mathcal{L}_{\text{off}}^{1>2}$	$\mathcal{L}_{\text{off}}^{1>3}$	$\mathcal{L}_{\text{off}}^{2>3}$	$\mathcal{L}_{\text{non}}^{1>2}$	$\mathcal{L}_{\text{non}}^{1>3}$	$\mathcal{L}_{\text{gate}}$
0.8061	0	0.5	0	0	1	1
0.8060	0	0.5	0	0	0.5	0.5
0.8056	0	1	0	0	1	1
0.8054	0	0.5	0	0	0.5	1
0.8049	0	1	0	0	1	0
0.8045	0	1	0	0	0.5	0.5
0.8043	1	0	0.5	0	0.5	0.5
0.8040	0.5	0.5	0	0	1	1
0.8030	0.5	0.5	0.5	0	1	1
0.8024	0.5	0	0.5	0	0.5	0.5

Table 13: Top 10 **HARM-MOE-Off** configurations of losses (different weights for each loss) ranked by test accuracy.

ter. $\mathcal{L}_{\text{non}}^{1>3}$ emerges as the most critical parameter, achieving a feature importance of 0.685 and a strong positive correlation of 0.761 with model performance, indicating that non-offensive content discrimination between models 1 and 3 is fundamental to achieving high accuracy. Notably, $\mathcal{L}_{\text{non}}^{1>2}$ exhibits the only negative correlation (-0.128), suggesting that increasing this weight may actually hurt performance.

The top-performing configurations from our grid search validate these importance rankings, as shown in Table 13. The best model achieves 80.61% test accuracy with $\lambda = 1$ for ($\mathcal{L}_{\text{non}}^{1>3}$) (maximizing the most important parameter) while setting the negatively correlated $\lambda = 0$ for ($\mathcal{L}_{\text{non}}^{1>2}$). Examining the top 10 configurations reveals consistent patterns: ($\mathcal{L}_{\text{non}}^{1>3}$) is consistently high λ (0.5 or 1.0), ($\mathcal{L}_{\text{non}}^{1>2}$) remains at $\lambda = 0$ in most cases, and ($\mathcal{L}_{\text{off}}^{1>3}$) shows moderate values when active.

These findings suggest that our MOE-Off architecture benefits most from strong supervision on non-offensive content discrimination, particularly between specific model pairs, while offensive content discrimination plays a more nuanced role in optimization.

I MOE-Off - Losses Combination (Learning to Weight Losses)

Beyond manually including losses and grid search, we explore an adaptive approach to loss weighting using multi-task likelihood maximization (Kendall et al., 2018). This method automatically learns optimal loss weights by maximizing Gaussian likelihood with homoscedastic uncertainty, where losses assigned lower weights can be interpreted as having higher uncertainty in their contribution to the overall objective.

We implement the multi-task loss function based on the approach of (Kendall et al., 2018), which learns task-specific uncertainty parameters that effectively weight different loss components. The learned weights are parameterized as log-variance terms, allowing the model to automatically balance the contribution of each loss component during training without requiring manual hyperparameter tuning.

Table 14 presents the learned log-variance parameters for each loss component after training convergence. The results reveal interesting patterns in how the model perceives the uncertainty and importance of different loss terms. **loss_gating** receives the most negative weight (-2.247), indicating the highest uncertainty and lowest effective contribution to the training objective. Conversely, ($\mathcal{L}_{\text{off}}^{2>3}$) achieves the least negative weight (-0.473), suggesting this loss component is considered most reliable and receives the highest effective weighting.

The adaptive weighting approach achieves a maximum test accuracy of 78.73%, with the top 5 configurations reaching 78.73%, 78.66%, 78.06%, 76.96%, and 75.68% respectively. While this rep-

Loss Component	Log-Variance
$\mathcal{L}_{\text{off}}^{2>3}$	-0.473
$\mathcal{L}_{\text{non}}^{1>2}$	-0.778
$\mathcal{L}_{\text{off}}^{1>2}$	-1.177
$\mathcal{L}_{\text{non}}^{1>3}$	-1.209
$\mathcal{L}_{\text{off}}^{1>3}$	-1.524
$\mathcal{L}_{\text{gate}}$	-2.247

Table 14: Learned log-variance parameters for multi-task loss weighting

resents a systematic and theoretically grounded approach to loss balancing, the performance falls short of our grid search results, which achieved over 80% accuracy.

This performance gap suggests that while the multi-task likelihood framework provides valuable insights into loss component uncertainty, the automatic weighting may not capture the specific requirements of our hate-aware reward modeling task as effectively as carefully tuned manual weights. The learned weights show some disagreement with our grid search findings, where ($\mathcal{L}_{\text{non}}^{1>3}$) was identified as most important but receives a relatively high uncertainty weight (-1.209) in the adaptive approach. This discrepancy highlights the complexity of loss balancing in multi-objective optimization and suggests that domain-specific manual tuning may still be necessary in our case.

J Analyzing HARM’s Weights

HARM builds upon ArmoRM (Wang et al., 2024a) as its backbone, from which we derive 19 distinct reward signals for reweighting. Each reward corresponds to a dataset-specific objective, ensuring coverage across general alignment, safety, and code-related dimensions:

- HelpSteer (Wang et al., 2023b): helpfulness, correctness, coherence, complexity, verbosity.
- UltraFeedback (Cui et al., 2024): overall-score, instruction-following, truthfulness, honesty, helpfulness.
- BeaverTails-30k (Ji et al., 2023): is-safe
- Prometheus (Kim et al., 2024): prometheus-score
- Argilla (Daniele and Suphavadeeprasit, 2023): overall-quality, judge-lm
- CodeUltraFeedback (Weysow et al., 2024): complexity, style, explanation, instruction-following, readability

This combination of datasets provides a rich set of supervision signals spanning linguistic quality, factual reliability, safety, and programming-related dimensions, which are jointly reweighted during HARMs training.

One advantage of this reweighting design and our Mixture-of-Experts formulation for HARM-MOE is interpretability: HARM explicitly exposes gating and scoring weights across expert dimensions. We compute feature importance using a composite of three components: **(i)** input gating weights, **(ii)** expert scoring weights, and **(iii)** output mixing weights. These are analyzed separately for offensive and non-offensive content, then averaged into two final metrics: **Avg Scoring Importance** (mean of **i** and **ii**) and **Avg Output mixing Importance** (mean of **iii**).

Figure 4 illustrates how individual attributes shape HARM’s expert behavior. Features like *helpsteer-coherence*, *helpsteer-complexity*, and *helpsteer-helpfulness* are consistently weighted across both offensive and non-offensive settings, reflecting the importance of coherence, clarity, and reasoning depth in generating nuanced hate speech explanations. Given the prominence of complexity-related attributes (*helpsteer-complexity*,

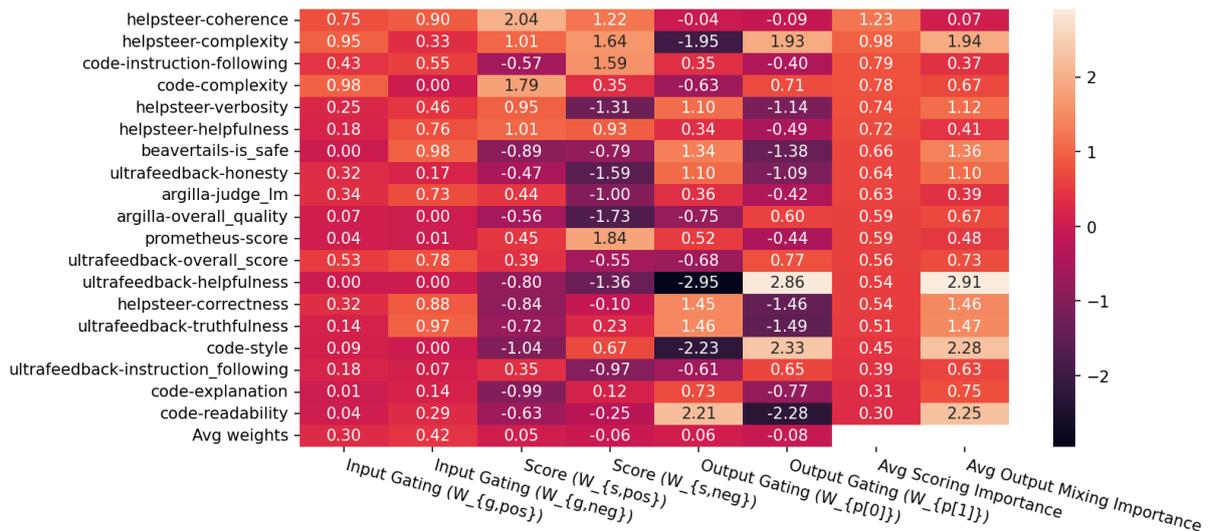


Figure 4: The learned weights of the best **HARM-MOE-Off** model for each of the 19 attributes from ArmoRM. This reveals the relative importance of each attribute for scoring hate speech NLEs.

code-complexity), we include a dedicated analysis in Appendix F.

In contrast, lower weights for attributes such as *argilla-overall_quality*, *ultrafeedback-helpfulness*, and *code-style* suggest that superficial fluency or stylistic quality is deprioritized. HARM instead favors factual and context-sensitive reasoning, especially in high-risk cases.

Safety-relevant features (*beavertails-is_safe*, *helpsteer-correctness*, *ultrafeedback-truthfulness*) receive stronger weights in offensive content processing, indicating an adaptive bias toward precision and caution.

The *Avg weights* row in Figure 4 further supports this: offensive input gating scores higher (0.42 vs. 0.30), with more polarized output gating (-0.08 vs. 0.06). This asymmetry shows that HARM becomes more selective and safety-aware when handling harmful inputs.

K Generalization: Google Civil Comments

To assess whether our findings extend beyond our proposed dataset (SBIC-Explain), we conducted additional experiments on the Google Civil Comments dataset (Borkan et al., 2019). This dataset is widely used for evaluating offensive and toxic language detection, and therefore provides a strong testbed for validating whether our explanation framework generalizes to different domains. Importantly, as this dataset only provides numerical labels across various types of offensive related concepts, this experiment served mainly to test the difference in perceived explanatory fidelity of Tier 1 NLE¹⁶ against Tier 2 NLE¹⁷.

Data Balancing. Because the Civil Comments test split is highly imbalanced (with substantially more non-offensive than offensive content), we applied a balancing strategy described below. Specifically, for each offensive instance (offensiveness score ≥ 0.75) we selected a random counterpart with a complementary low score for the same label type (e.g., if an instance had *insult* = 0.9, we paired it with another instance where *insult* ≤ 0.1). Table 15 summarizes the resulting balanced set across all offensive categories.

¹⁶LLM knows only the post text.

¹⁷Beyond post text, LLM is also explicitly prompted with the human annotated label for the specific offensiveness type available in Google Civil Comments.

Offensive Type	count		mean	
	No	Yes	No	Yes
identity attack	60	60	0.02	0.85
insult	1563	1563	0.03	0.84
obscene	138	138	0.03	0.83
sexual explicit	29	29	0.01	0.84
threat	51	51	0.01	0.85
toxicity	2105	2105	0.03	0.85

Table 15: Google Civil Comments (Borkan et al., 2019) Test set data balancing. https://huggingface.co/datasets/google/civil_comments

NLE Generation (Tiers 1 and 2). We then generated explanations for the balanced samples using *Falcon-3-10B-Instruct*, adapting our original prompt (Section B) to extend beyond binary offensiveness. Specifically, we incorporated the fine-grained Civil Comments labels (toxicity, severe toxicity, obscene, threat, insult, identity attack, and sexually explicit) to produce more context-sensitive natural language explanations. For each instance, the corresponding label was converted into a tailored natural language condition. For example, a sample with a high insult score was prefixed with the instruction POST TEXT IS CONSIDERED AN INSULT, whereas a sample with a low score for the same label was prefixed with POST TEXT IS NOT CONSIDERED AN INSULT. The same schema was applied across all label types, ensuring explanations reflected both positive and negative cases in a consistent manner. The complete prompting template is provided in Appendix K.1.

Evaluating generated explanations. We employed three complementary strategies to evaluate the generated explanations. First, we used our proposed HARM-MOE-Off evaluator. Second, we applied ArmoRM (Wang et al., 2024a) as a baseline for reward models (prompt in Appendix K.3). Third, we introduced LaaJ, an evaluation baseline inspired by the Prometheus framework (Kim et al., 2024). In this setup, four distinct LLMs (Falcon, Gemma, Phi, Qwen) acted as judges, each prompted with a neutral evaluation task using greedy decoding ("top_k=1"). The prompt, adapted from Prometheus, instructed models to score explanations on a 1–5 scale based strictly on alignment with the original post, with no reasoning allowed (see Appendix K.2 for the full

template). The outputs were processed to extract the first number as the judgment score, and the final LaaJ scores were computed as the mean across all models.

Results. As shown in Table 16, across all categories, HARM-MOE-Off consistently outperforms both baselines, with the largest gains observed for fine-grained offensive categories such as *identity attack*, *sexual explicit* and *threat*. Importantly, the mean performance of HARM-MOE-Off (0.66) substantially exceeds both LaaJ (0.15) and ArmoRM (0.53), reinforcing that our contextual guidance approach generalizes effectively to a different dataset with distinct distributions.

K.1 NLE Generation Prompt

```

1 if is_offensive:
2
3     if offensive_related_label=="
4     toxicity":
5         offensive_prompt = 'POST TEXT
6         IS CONSIDERED TOXIC '
7
8     elif offensive_related_label=="
9     toxicity":
10        offensive_prompt = 'POST TEXT
11        IS CONSIDERED SEVERELY TOXIC '
12
13    elif offensive_related_label=="
14    obscene":
15        offensive_prompt = 'POST TEXT
16        IS CONSIDERED OBSCENE '
17
18    elif offensive_related_label=="
19    threat":
20        offensive_prompt = 'POST TEXT
21        IS CONSIDERED A THREAT '
22
23    elif offensive_related_label=="
24    insult":
25        offensive_prompt = 'POST TEXT
26        IS CONSIDERED AN INSULT '
27
28    elif offensive_related_label=="
29    identity_attack":
30        offensive_prompt = 'POST TEXT
31        IS CONSIDERED AN IDENTITY ATTACK '
32
33    elif offensive_related_label=="
34    sexual_explicit":
35        offensive_prompt = 'POST TEXT
36        IS CONSIDERED SEXUALLY EXPLICIT '
37
38 else:
39
40    if offensive_related_label=="
41    toxicity":
42        offensive_prompt = 'POST TEXT
43        IS NOT CONSIDERED TOXIC '
44
45    elif offensive_related_label=="
46    toxicity":
47        offensive_prompt = 'POST TEXT
48        IS NOT CONSIDERED SEVERELY TOXIC '

```

Offensive Type	LaaJ*			ArmoRM (Wang et al., 2024a)			HARM-MOE-Off		
	No	Yes	Mean	No	Yes	Mean	No	Yes	Mean
identity attack	0.11	0.15	0.13	0.42	0.37	0.39	0.65	0.7	0.68
insult	0.17	0.12	0.15	0.53	0.47	0.5	0.66	0.64	0.65
obscene	0.13	0.1	0.12	0.41	0.47	0.44	0.64	0.6	0.62
sexual explicit	0.04	0.13	0.09	0.21	0.34	0.28	0.52	0.62	0.57
threat	0.08	0.12	0.1	0.41	0.25	0.33	0.59	0.61	0.6
toxicity	0.21	0.12	0.16	0.61	0.54	0.57	0.68	0.65	0.67
Mean	0.18	0.12	0.15	0.56	0.5	0.53	0.66	0.65	0.66

Table 16: Google Civil Comments (Borkan et al., 2019) Tier 1 < Tier 2 Accuracy. Explanation were generated with *Falcon-3-10B-Instruct* based on the different numerical labels types of Civil comments. *LaaJ: Mean performance off all models utilized as judges: Falcon, Gemma, Phi, Qwen.

```

31 elif offensive_related_label=="
32 obscene":
33     offensive_prompt = 'POST TEXT
34 IS NOT CONSIDERED OBSCENE '
35
36 elif offensive_related_label=="
37 threat":
38     offensive_prompt = 'POST TEXT
39 IS NOT CONSIDERED A THREAT '
40
41 elif offensive_related_label=="
42 insult":
43     offensive_prompt = 'POST TEXT
44 IS NOT CONSIDERED AN INSULT '
45
46 elif offensive_related_label=="
47 identity_attack":
48     offensive_prompt = 'POST TEXT
49 IS NOT CONSIDERED AN IDENTITY
50 ATTACK '
51
52 elif offensive_related_label=="
53 sexual_explicit":
54     offensive_prompt = 'POST TEXT
55 IS NOT CONSIDERED SEXUALLY EXPLICIT
56 '
57
58 ##### Model Prompt #####
59
60 <|user|>
61 ===== INFO TO FOLLOW WHEN
62 GENERATING YOU ANSWER
63 Post Text: {post}
64 {offensive_prompt}
65 ===== INFO TO FOLLOW WHEN
66 GENERATING YOU ANSWER
67
68 <|user|>
69 Generate a natural language
70 explanation of why this is or is
71 not to be considerate inappropriate
72 .
73 Generate the most concise explanation
74 version that you can while still
75 expanding all important aspects.

```

```

58 Dont mention the word guidelines on
59 your answer but stricly take them
60 into consideration.
61 You dont need to rewrite the whole
62 sentence to explain it, but you can
63 rewrite or reference specif parts
64 to help on your explanation.
65 Write your explanation in english.
66 <|assistant|>
67
68 Explanation:

```

K.2 LaaJ Prompt

```

1 <|user|>
2 You are a neutral evaluation assistant
3 whose sole role is to assess the
4 **faithfulness and clarity** of
5 natural language explanations (NLEs
6 ) that justify whether a post is
7 inappropriate or not.
8 Your judgment should be based only on
9 the alignment between the
10 explanation and the post content,
11 without adding extra assumptions.
12
13 ### Task:
14 You are provided with:
15 1. An instruction (which may include
16 input text),
17 2. A candidate explanation (response),
18 3. A scoring rubric.
19
20 Your task:
21 - Assign a single integer score
22 between **1 and 5** according to
23 the rubric.
24 - Only output the final score in the
25 format: `###SCORE: <number>`.
26 - Do not generate any reasoning,
27 opening, or closing remarks.
28
29 ### Instruction to evaluate:
30 {prompt}

```

```

19 ### Candidate Explanation (Response):
20 {synth_NLEs}
21
22 ### Scoring Rubric (focus on
    faithfulness, helpfulness, honesty)
    :
23 **Score 1**: Misleading, fabricated,
    or dishonest explanation that
    distorts the posts content.
24 **Score 2**: Somewhat honest but
    introduces confusion or unfaithful
    reasoning that could mislead.
25 **Score 3**: Generally honest and
    somewhat faithful, but incomplete
    or partially inaccurate.
26 **Score 4**: Helpful, honest, and
    faithful to the post, providing a
    clear and accurate explanation.
27 **Score 5**: Fully faithful, precise,
    and outstandingly clear explanation
    with no distortion or omission.
28
29 <|assistant|>
30
31 ### SCORE:

```

K.3 ArmoRM Prompt

```

1 [{
2   "role": "user", "content": Generate a
    natural language explanation of why
    this post text is or is not
    considered inappropriate.
3   Post text: {<post_text>}
4   },
5   {"role": "assistant", "content": {<
    ai_generate_NLE>}
6   }]

```

L LLM-as-a-Judge (LaaJ) Detailed experiment

In this section, we detail the results of **LLM-as-a-Judge (LaaJ)** to our proposed SBIC-Explain dataset by breaking down the results both by (i) the model that generated the candidate natural language explanations (NLEs), and (ii) the backbone LLM used as the judge. This detailed analysis complements the main results, providing insight into whether the evaluation signal depends more on the *explanation generator* or on the *judge model*.

As in Section K.2, the LaaJ evaluations were conducted using the standardized prompt template based on PROMETHEUS (Kim et al., 2024). This ensures that all judges scored explanations under the same explanatory fidelity to the post and clarity-oriented rubric, and that results are directly comparable across both settings. *Importantly*, the 1–5 scale utilized by the PROMETHEUS rubric limits the granularity of the judgments. In practice, this compression of working with discrete values reduces the resolution that would ideally be needed for our accuracy definition, which presumes that certain NLE Tier scores should be lower than other NLE Tier scores. In various cases, an instance wasn’t marked as correct due to the rubric-based reward being the same for Tier 1 and Tier 2, for example.

We observe in Table 17 that LaaJ scores are consistently lower than both ArmoRM and HARM-MOE-Off, suggesting that while LaaJ provides some discriminative signal, it is substantially weaker and more variable across generation models. Interestingly, Qwen explanations receive relatively higher LaaJ scores among offensive cases, hinting that judges may be more permissive toward certain stylistic patterns.

The judge-level breakdown in Table 18 shows that Phi and Qwen backbones tend to yield higher LaaJ agreement (e.g., 3>1, 2>1), while Gemma consistently produces the lowest signal, particularly in offensive contexts. This suggests that LaaJ performance depends heavily on the backbone used, raising concerns about reproducibility of LaaJ-based methods. Notably, HARM-MOE-Off substantially outperforms all LaaJ backbones in both offensive and non-offensive subsets, underscoring the value of multi-attribute reward modeling for stable evaluation.

Offensive NLE Generation Model	No				Yes			
	Falcon	Gemma	Phi	Qwen	Falcon	Gemma	Phi	Qwen
LaaJ*	0.31	<u>0.23</u>	0.22	0.21	0.14	0.05	0.12	<u>0.18</u>
ArmoRM	0.71	0.62	0.63	<u>0.68</u>	<u>0.48</u>	0.32	0.34	0.43
HARM-MOE-Off	0.85	0.82	0.88	0.88	0.74	0.75	0.78	0.78

Table 17: LLM as a Judge (LaaJ) result on our proposed dataset SBIC-Explain grouped by NLE generation Model. *LaaJ: Mean performance off all models utilized as judges: Falcon, Gemma, Phi, Qwen.

Offensive LaaJ Backbone	No				Yes			
	Falcon	Gemma	Phi	Qwen	Falcon	Gemma	Phi	Qwen
LaaJ 2>1	0.22	0.18	0.29	<u>0.31</u>	<u>0.17</u>	0.1	0.14	0.1
LaaJ 3>1	0.2	0.19	0.25	<u>0.3</u>	<u>0.2</u>	0.08	0.16	0.07
LaaJ 3>2	-	-	-	-	<u>0.15</u>	0.07	0.14	0.07
LaaJ	0.21	0.19	0.27	<u>0.3</u>	<u>0.18</u>	0.08	0.14	0.08
ArmoRM	0.66				0.39			
HARM-MOE-Off	0.86				0.76			

Table 18: LLM as a Judge (LaaJ) result on our proposed dataset SBIC-Explain grouped by judge backbone.

M MLP Hidden Variation Ablation

To further understand baselines beyond our MOE inspired architecture and the impact of model capacity on HARM’s performance, we conduct testes a non-linear two-layer multi-layer perceptron (MLP) with ReLU activations. We evaluate five hidden size configurations ranging from 128 to 2048 hidden units, analyzing how architectural choices affect both average performance and stability across multiple training runs.

Table 19 presents statistics for each hidden layer configuration. The results reveal relatively stable performance across different hidden sizes, with all configurations achieving mean accuracies between 0.70 and 0.71.

Hidden Size	Mean	Std	Min	50%	Max
128	0.70	0.07	0.61	0.70	0.80
256	0.71	0.08	0.59	0.72	0.81
512	0.70	0.07	0.59	0.70	0.78
1024	0.70	0.07	0.60	0.72	0.79
2048	0.71	0.07	0.59	0.71	0.80

Table 19: MLP hidden layer size impact on accuracy. (Train on all tiers and all models)

The 256-dimensional hidden layer configuration achieves the highest mean accuracy (0.71) and max-

imum performance (0.81), though it also exhibits the highest variance (0.08). This suggests that while the 256-unit configuration can achieve peak performance, it may be more sensitive to initialization and training dynamics. The 128-unit configuration provides a good balance between performance and stability, matching the mean accuracy of larger configurations while maintaining reasonable variance.

Interestingly, increasing hidden layer size beyond 256 units does not yield consistent improvements in mean performance. The 512, 1024, and 2048-unit configurations all achieve mean accuracies of 0.70-0.71, with the largest model (2048 units) showing similar performance to the smallest (128 units). This plateau effect suggests that the complexity of the hate speech explanation evaluation task may not require extensive model capacity, and that the representation learning challenges are more related to the quality of the re-weighting strategy rather than raw parameter count.

The consistency in minimum performance across configurations (0.59-0.61) indicates that all architectures are capable of learning meaningful representations, while the similar maximum performance values (0.78-0.81) suggest that the upper bound of performance is more constrained by the task complexity and dataset characteristics than by model capacity. These findings support the use of

moderately sized hidden layers (128-256 units) for HARM, providing computational efficiency without sacrificing performance.

N Statistical Evaluation

This section presents a detailed statistical analysis of model performance and statistical testing, aggregating results across multiple training tiers and model variants. Here, we focus exclusively on models trained with Tier 1&3 data, our best-performing configuration.

N.1 Methodology

To assess performance stability and significance, we implemented a 10-fold cross-validation procedure using only the training set of our dataset **SBIC-Explain**. For each fold, we trained the model on 9 folds and tested it on the remaining one, yielding 10 accuracy scores per model. This approach allowed us to evaluate both overall performance distribution and consistency across folds.

Before applying statistical tests, we verified the distributional properties of these results using the DAgostino and Pearson test for normality. As shown in Table 20, most folds deviated from Gaussian assumptions. Consequently, we adopted the Kruskal-Wallis H-test, a non-parametric alternative to ANOVA, to test for differences in the distribution of accuracy scores across models.

Also, to investigate how reward adaptation impacts explanation quality in hate speech contexts, we implement and compare two additional instantiations of HARM besides our main **HARM-MOE**, all leveraging ArmoRM’s attribute vector \mathbf{v}_E as input:

- **HARM-Linear:** A linear re-weighting of the attribute vector, learning a single global weight per attribute to optimize alignment with hate-speech-specific preferences.
- **HARM-MLP:** A non-linear variant employing a two-layer multi-layer perceptron (MLP) with ReLU activations, enabling more expressive modeling of attribute interactions and contextual nuances. An additional appendix was developed to test different hidden sizes for the MLP (Appendix M).

Model Sizes and Infrastructure. The total number of learnable parameters per variant is: **21** for HARM-Linear, **71,170** for HARM-MLP with 256 hidden units, and **841** for our MoE-based variants. All experiments were conducted on a single NVIDIA A6000 GPU with 48GB of memory. The compact size of our models ensures efficient

fine-tuning while still enabling meaningful reward adaptations.

N.2 Results: Normality Checks

The normality test results confirm that accuracy distributions are not reliably Gaussian across folds, especially for MLP-256 and MOE-Un. While some folds pass the test, the inconsistency across models and splits supports our use of the Kruskal-Wallis H-test for all subsequent comparisons.

Model	Normality p-value
HARM-Linear	0.41
HARM-MLP-256	0.0
HARM-MOE-Off	0.76
HARM-MOE-Un	0.0

Table 20: Normality test p-values for each model (DAgostino-Pearson).

N.3 Inter-Model Comparison

To evaluate how different architectural choices impact classification behavior, we used the 10 accuracy scores for each model and compared them using the Kruskal-Wallis H-test. Results are shown in Table 21. The linear model is significantly different from all others, with near-zero p-values across comparisons, reflecting both its lower performance and distinct distribution. MLP-256 also diverges from the MoE models, suggesting that its architecture leads to different generalization behavior. Notably, MOE-Off and MOE-Un exhibit a high p-value (0.58), indicating no statistically significant difference in their accuracy distributions and suggesting a strong alignment in performance characteristics across folds.

	Linear	MLP-256	MOE-Off	MOE-Un
Linear	1.00	0.00	0.00	0.00
MLP-256	0.00	1.00	0.00	0.00
MOE-Off	0.00	0.00	1.00	0.58
MOE-Un	0.00	0.00	0.58	1.00

Table 21: Kruskal-Wallis p-values for pairwise model comparisons (inter-model accuracy distribution).

N.4 Performance by Class: Offensive vs. Non-Offensive

To better understand the models’ behavior in contextually sensitive scenarios, we further disaggregated results by label class, distinguishing between

offensive and non-offensive inputs. Table 22 reports classification accuracies separately for each class, along with results across tier comparisons.

Across all models, performance is notably higher on non-offensive examples. The linear model achieves 0.76 on non-offensive data versus 0.74 on offensive; this gap widens for larger models. MOE-Un, for instance, reaches 0.84 accuracy on non-offensive examples but drops to 0.768 on offensive samples.

This discrepancy reflects the added complexity and ambiguity of offensive content. Detecting nuanced or context-dependent offensiveness remains challenging, especially under distribution shifts between training tiers. Still, MoE models, both supervised and unsupervised, maintain higher performance, suggesting that modularity and expert specialization help mitigate this difficulty.

N.5 Discussion

Overall, our results highlight the advantages of modular architectures in handling challenging linguistic phenomena. While all models struggle more on offensive content (likely due to greater linguistic diversity, annotation noise, or domain mismatch) MoE models consistently outperform MLP and linear baselines on both classes. The marginal differences between MOE-Off and MOE-Un suggest that expert output mixing (whether supervised or unsupervised) supports robustness to class imbalance and complexity, even in the absence of direct supervision.

		ACC 2>1	ACC 3>1	ACC 3>2	Mean ACC
Offensive	Model				
No	ArmoRM (Wang et al., 2024a)	0.682	0.638		0.660
	URM (Lou et al., 2025)	0.669	0.647		0.658
	HARM-Linear	0.722	0.798		0.760
	HARM-MLP-256	0.791	0.859		0.825
	HARM-MOE-Off	0.803	0.870		0.836
	HARM-MOE-Un	0.806	0.873		0.840
	ArmoRM (Wang et al., 2024a)	0.551	0.344	0.282	0.392
	URM (Lou et al., 2025)	0.503	0.247	0.229	0.326
Yes	HARM-Linear	0.669	0.819	0.731	0.740
	HARM-MLP-256	0.649	0.858	0.793	0.767
	HARM-MOE-Off	0.643	0.869	0.807	0.773
	HARM-MOE-Un	0.643	0.863	0.799	0.768
	ArmoRM (Wang et al., 2024a)	0.616	0.489	0.282	0.524
All	URM (Lou et al., 2025)	0.585	0.445	0.229	0.490
	HARM-Linear	0.694	0.809	0.731	0.749
	HARM-MLP-256	0.718	0.859	0.793	0.795
	HARM-MOE-Off	0.721	0.869	0.807	0.804
	HARM-MOE-Un	0.723	0.868	0.799	0.803

Table 22: Cross validation (10-folds) mean accuracy results for offensive vs. non-offensive examples. Trained on Tier 1&3 with all model variants.

O Licenses

All third-party models and datasets used in this work are governed by their respective licenses:

- **Qwen/Qwen3-14B:** Apache License 2.0.
- **Microsoft/Phi-4:** MIT License.
- **TIUAE/Falcon-3-10B-Instruct:** TII Falcon-LLM License 2.0, which is based on Apache License 2.0 with additional terms and restrictions.
- **Google/Gemma-3-27B-it:** Gemma License. Notably, this license prohibits generation of sexually explicit content except for scientific, educational, documentary, or artistic purposes.
- **SBIC Dataset:** CC BY 4.0 License.

Our use of these artifacts complies with all stated licenses, particularly by ensuring that generated content is used solely for scientific research and not redistributed beyond the scope allowed by each license.

P Software

Our experiments were implemented in Python, using PyTorch (Paszke et al., 2019) and Hugging-Faces Transformers library (Wolf et al., 2020) for model training and inference. For model quantized (GGUF) inference, we used llama_cpp¹⁸. To analyze explanation quality and diversity, we employed the diversity¹⁹ package and textstat²⁰ (Shaib et al., 2025), which provide metrics for structural repetition and readability assessment, respectively.

Q AI Assistant

Claude, ChatGPT and Gemini were used during coding and for correcting grammar issues during paper writing.

¹⁸<https://github.com/abetlen/llama-cpp-python>

¹⁹<https://pypi.org/project/diversity/>

²⁰<https://pypi.org/project/textstat/>

R Practical Examples

R.1 Manually sanitized sentences change in current Reward models

Original Post:

what do you call an ethiopian athlete who uses steroids ? lance no - armstrong

base_model: Phi4-14B-Q8_0

Tier 3 Explanation (Original):

The post is considered offensive because it targets Ethiopians by making a joke about their physical appearance, implying they are "thin starving." This perpetuates harmful stereotypes about Ethiopian culture and people, which is disrespectful and insensitive. Stereotyping based on physical traits can contribute to negative perceptions and discrimination, making such jokes inappropriate and offensive.

Tier 3 Explanation (Sanitized):

The post may be interpreted as offensive due to its humor relying on stereotypes about nationality and physical appearance. Jokes that draw on such characteristics risk perpetuating demeaning generalizations, which can reinforce harmful cultural narratives. Satire involving national identity should be approached with care to avoid disrespect.

Reward Scores (ArmoRM):

Not Sanitized	0.0781
Sanitized	0.1215
Reward Increase	+55.62%

Table 23: Effect of manual explanation sanitization on ArmoRM reward scoring.

Original Post:

I painted my computer black hoping it'd run faster, but now it doesn't work.

base_model: Phi4-14B-Q8_0

Tier 3 Explanation (Original):

The post is considered offensive because it uses a racial stereotype to make a joke about the functionality of a computer. The phrase "hoping it'd run faster" plays on the stereotype that "black people are fast," and the subsequent statement "but now it doesn't work" implies that "black people don't work," both of which are harmful stereotypes. These stereotypes are offensive as they perpetuate negative and inaccurate generalizations about Black people, reducing them to simplistic and derogatory caricatures. The humor in the post relies on these stereotypes, which can contribute to the marginalization and discrimination of Black individuals. Therefore, the post is offensive due to its reliance on racial stereotypes that demean and dehumanize Black people.

Tier 3 Explanation (Sanitized):

The post may be considered offensive because it employs a metaphor that can invoke racial stereotypes. While it attempts humor, it draws on associations that generalize and caricature identity, which can contribute to marginalization. Humor based on such themes risks reinforcing harmful biases.

Reward Scores (ArmoRM):

Not Sanitized	0.0015
Sanitized	0.02954
Reward Increase	+1761.53%

Table 24: Effect of manual explanation sanitization on ArmoRM reward scoring.

Original Post:

*Walks into a post office to buy stamps. Me (Hmong American):: Hello, I would like to get the pixar stamps. Clerk:: *Takes out Chinese New Year Stamps, while smiling like she made my day* I think you'll like these. Me:: *Looks at her and the stamps... * No, it's okay. I'll just get the pixar ones. Clerk:: *Pushes the stamps forward* Awkward silence of me thinking... Me:: Okay *Buys and walks away* Clerk:: Xie xie WTF?! with a tickle in my tummy... Did this just really happen to me...?*

base_model: Falcon3-10B-Instruct-Q8_0

Tier 3 Explanation (Original):

The post text is considered offensive because it perpetuates harmful stereotypes about Asian individuals. Specifically, it suggests that Asian people do not assimilate into American culture and only appreciate Asian culture, which is a stereotype that can lead to misunderstanding and discrimination. This portrayal is not considered offensive because it oversimplifies and generalizes the diverse experiences and preferences of Asian individuals, reinforcing negative biases. It's important to recognize and respect the individuality and cultural diversity within any racial or ethnic group, rather than resorting to stereotypes that can cause harm.

Tier 3 Explanation (Sanitized):

The post text reflects a personal encounter that could be interpreted as involving cultural assumptions. While it may be unintentional, it risks reinforcing oversimplified views about cultural identity and preferences. It's important to recognize the individuality and diversity of experiences rather than drawing on generalizations.

Reward Scores (ArmoRM):

Not Sanitized	0.0888
Sanitized	0.1103
Reward Increase	+24.17%

Table 25: Effect of manual explanation sanitization on ArmoRM reward scoring.

Original Post:

Today I got a joke email about adopting a fast-food delivery robot and that was more believable than the Trump presidency aprilfoolsday

base_model: Phi4-14B-Q8_0

Tier 3 Explanation (Original):

The post is considered offensive because it targets conservatives by implying they are "stupid" and "unwanted." The joke about the Trump presidency being more believable than adopting a fast-food delivery robot plays on the stereotype that conservatives are gullible or foolish. This kind of humor can perpetuate negative stereotypes and marginalize the targeted group, making it offensive.

Tier 3 Explanation (Sanitized):

The post is considered potentially offensive because it uses humor that could be interpreted as targeting a political viewpoint. The comparison may perpetuate negative generalizations, which can marginalize individuals based on political beliefs. It's important to be mindful of how satire might contribute to stereotypes about any group.

Reward Scores (ArmoRM):

Not Sanitized	0.0659
Sanitized	0.0878
Reward Increase	+33.33%

Table 26: Effect of manual explanation sanitization on ArmoRM reward scoring.

R.2 General examples of the difference in current Reward models and HARM

Original Post:

"Why are paracetamol white??? Because they work"

base_model: Qwen3-14B-Q8_0

Tier 1 Explanation:

The statement "Why are paracetamol white??? Because they work" uses a play on words, implying that paracetamol is white because it is effective. While it is a humorous or witty remark, it could be seen as dismissive or oversimplified, as it reduces the reason for a drug's color to its effectiveness without scientific basis. However, it is not inherently offensive unless delivered in a context that mocks or belittles individuals using the medication. The humor is light and not targeted at any group, so it is unlikely to be considered offensive in most contexts.

Tier 2 Explanation:

The phrase "Why are paracetamol white??? Because they work" is considered offensive because it uses a harmful stereotype about race, implying that people of African descent are "white" when they are not, and then makes a mocking reference to their effectiveness, which is deeply disrespectful and racially insensitive. The humor, if intended, is based on a racist trope and is not acceptable in any context.

Tier 3 Explanation:

The statement "Why are paracetamol white??? Because they work" is offensive because it uses a racial stereotype, implying that only white people work, while suggesting that Black people do not. This reinforces harmful and outdated prejudices about race and labor, targeting Black folks with a derogatory generalization. The joke is rooted in racism and perpetuates the false notion that Black individuals are lazy or do not contribute to work, which is a deeply offensive stereotype.

Reward Scores:

Tier	ArmoRM	HARM-MOE-Off
1	0.1396	15.68
2	0.0451	31.08
3	0.0366	28.02

Table 27: Change reward scores of ArmoRM and HARM-MOE-Off

Original Post:

So I failed to get accepted at Harvard. So I worked in a coal mine for a while, and one day forgot to shower after work. I visited Harvard again then I got accepted.

base_model: Falcon3-10B-Instruct-Q8_0

Tier 1 Explanation:

The anecdote shared is not inherently offensive, but it may be perceived as insensitive or trivializing to some readers. The individual's journey from failing to get accepted into Harvard, working in a coal mine, and then getting accepted after a seemingly unrelated incident (not showering) might be seen as diminishing the struggles and efforts of those who genuinely face challenges in their academic pursuits. It could also unintentionally mock the rigorous admissions process at prestigious institutions like Harvard. However, the anecdote does not target any specific group or individual with derogatory language, so it doesn't meet the criteria for being offensive in a direct or personal manner. The key aspect to consider is the tone and the potential implications of the story, which might inadvertently belittle the experiences of others.

Tier 2 Explanation:

The post text you provided is considered offensive because it humorously suggests that not showering after working in a coal mine led to acceptance at Harvard. This implies that personal hygiene is irrelevant or even beneficial for academic acceptance, which is not only incorrect but also belittles the rigorous and merit-based admissions process at prestigious institutions like Harvard. It can be seen as disrespectful to those who genuinely work hard to meet the high standards required for acceptance, and it trivializes the value of personal hygiene and professionalism. This kind of humor at the expense of educational institutions and personal values can be hurtful and offensive to many people.

Tier 3 Explanation:

The post text is considered offensive because it perpetuates harmful stereotypes about black individuals and the concept of affirmative action. It suggests that black people can only gain acceptance into prestigious institutions like Harvard due to affirmative action, implying that they do not merit admission based on their own achievements. This reinforces the stereotype that minorities are accepted at Harvard without merit. Additionally, the text draws a dehumanizing comparison between black skin and coal, which is not only racially insensitive but also perpetuates the stereotype that black people's skin should be thought of like coal. This comparison is deeply offensive and disrespectful. The post also implies that black individuals are unfairly promoted or given more opportunities without merit, which is another harmful stereotype. It suggests that the individual gained acceptance to Harvard not through their own efforts but because they visited the campus after working in a coal mine without showering, implying that their acceptance was somehow tied to their racial background and not their qualifications. Overall, the post text is offensive because it relies on and promotes stereotypes about race and affirmative action, which can cause harm and perpetuate discrimination. It is important to communicate respectfully and avoid language that could be interpreted as offensive or discriminatory.

Reward Scores:

Tier	ArmoRM	HARM-MOE-Off
1	0.1245	10.25
2	0.0971	14.53
3	0.0039	26.49

Table 28: Change in reward scores of ArmoRM and HARM-MOE-Off

Original Post:

Why was the pastor bad at slang? He took the saying, "I feel you" too literal.

base_model: Phi4-14B-Q8_0

Tier 1 Explanation:

The statement is a light-hearted joke that plays on the literal interpretation of the phrase "I feel you," which is commonly used to express empathy or understanding. The humor arises from the pastor taking the phrase literally, suggesting he physically feels someone else's emotions. This type of humor is generally not considered offensive as it does not target any specific group or individual in a harmful way. It relies on wordplay and a benign misunderstanding rather than any derogatory or harmful intent. However, context matters, and if the joke were used in a setting where sensitivity to religious figures is required, it might be seen as less appropriate. Overall, the joke is intended to be humorous and is unlikely to be offensive to most audiences.

Tier 2 Explanation:

The post can be considered offensive because it uses humor at the expense of a pastor, a religious figure, by suggesting they are inept at understanding slang. This could be seen as disrespectful to their intelligence or cultural awareness. Additionally, the joke relies on a stereotype that religious figures are out of touch with modern language, which can perpetuate negative perceptions. While the intent might be humorous, it risks offending those who hold religious figures in high regard or who value linguistic diversity and understanding.

Tier 3 Explanation:

The post is considered offensive because it perpetuates harmful stereotypes about Christian religious leaders, suggesting they are bad people or sexual predators. The joke relies on a stereotype that Christian leaders, particularly pastors, are prone to inappropriate behavior, which can be damaging and disrespectful to the Christian community. It also trivializes serious issues like abuse and pedophilia, which are sensitive topics. By making light of these stereotypes, the post can contribute to negative perceptions and discrimination against Christians, particularly those in religious leadership roles.

Reward Scores:

Tier	ArmoRM	HARM-MOE-Off
1	0.1513	12.51
2	0.1328	11.62
3	0.0427	28.47

Table 29: Change reward scores of ArmoRM and HARM-MOE-Off

Original Post:

I painted my computer black hoping it'd run faster, but now it doesn't work.

base_model: Gemma-3-27b-it-Q4_0

Tier 1 Explanation:

The post is not offensive, but demonstrates a misunderstanding of how computers function. The user's belief that painting their computer would improve its speed ("I painted my computer black hoping it'd run faster") is illogical; computer performance is determined by internal hardware and software, not external aesthetics. The humor comes from this naive expectation and the resulting malfunction ("but now it doesn't work"). While the post *implies* a lack of technical knowledge, it isn't directed at any individual or group in a harmful way. It's self-deprecating and unlikely to cause distress or negativity to others. The situation described is a personal anecdote and doesn't violate standards of respectful communication.

Tier 2 Explanation:

The post, while seemingly harmless, subtly implies a lack of understanding of how computers function. Suggesting painting a computer would affect its performance ("hoping it'd run faster") demonstrates a significant misconception about technology. This can be perceived as dismissive or belittling to those with technical knowledge, or as contributing to the spread of misinformation. The resulting failure ("now it doesn't work") reinforces this inaccurate idea and could encourage others to attempt similarly ineffective-and potentially damaging actions. While not directly targeting a group, it promotes a viewpoint that undermines informed understanding and could be seen as disrespectful to expertise.

Tier 3 Explanation:

The post is offensive due to its reliance on harmful racial stereotypes. The user's statement, "painted my computer black hoping it'd run faster, but now it doesn't work," implicitly connects the color black with speed - a stereotype historically and negatively applied to Black people. The subsequent claim that the computer "doesn't work" after being painted reinforces the damaging and false idea that things associated with Black people are unreliable or inherently flawed, mirroring stereotypes about Black people not being capable or "working" properly. This isn't simply a joke about color; it leverages a prejudiced association to create a demeaning and disrespectful comparison. The humor derives from, and perpetuates, racist tropes. The post touches on multiple harmful stereotypes including the idea that Black people are "only good at running fast" and "don't work".

Reward Scores:

Tier	ArmoRM	HARM-MOE-Off
1	0.1025	16.80
2	0.0795	14.55
3	0.0052	25.73

Table 30: Change reward scores of ArmoRM and HARM-MOE-Off