

ClaHF: A Human Feedback-inspired Reinforcement Learning Framework for Improving Classification Tasks

Anonymous ACL submission

Abstract

Text classification models are typically trained via supervised fine-tuning (SFT). However, SFT essentially performs behavior cloning from instance-wise labels and thus fails to adequately capture relative preference relations among samples, which limits the model’s ability to shape decision boundaries and calibrate predictive confidence. In this paper, we propose ClaHF, a human feedback-inspired reinforcement learning (RL) framework for text classification that integrates preference modeling and RL optimization into the classification pipeline without requiring additional human annotations. Unlike prior work that relies solely on instance-wise supervision, ClaHF constructs multiple candidate predictions together with their relative ranking relations, and jointly models the Top-1 preference and the ordering among non-optimal candidates within a reward model (RM). This design converts conventional label supervision into preference signals that are directly applicable to policy optimization. We conduct systematic evaluations on eight classification tasks spanning three categories of scenarios. Results demonstrate that ClaHF consistently improves both classification performance and confidence calibration across diverse language models (LMs). The data and code are available at <https://anonymous.4open.science/r/ClaHF>.

1 Introduction

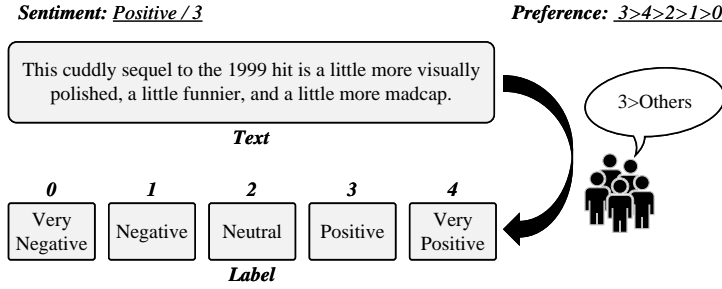
Text classification is one of the most fundamental and important tasks in natural language processing (NLP), forming the basis of a wide range of downstream applications such as sentiment analysis, news categorization and emotion recognition (Li et al., 2022). In recent years, the rapid development of pre-trained LMs has led to remarkable improvements in these domains, substantially boosting downstream performance (Sun et al., 2023). Nevertheless, most existing text classification mod-

els are still trained via SFT, where model predictions are aligned with gold labels on an instance-by-instance basis (Fonseca et al., 2025; Wei and Zhu, 2025). From a learning paradigm perspective, SFT essentially performs behavior cloning: the model is trained to imitate annotators’ decisions on the training data in order to learn classification boundaries (Chu et al., 2025).

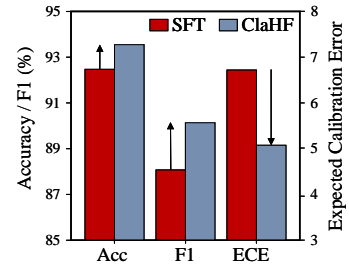
Although SFT effectively activates the semantic representation capacity of pre-trained models, it lacks further exploration of the decision space (Fu et al., 2023). When data quality is limited, when class boundaries differ only in subtle semantics, or when annotations themselves contain bias, SFT tends to amplify these issues. As a result, the model may overfit the label distribution rather than the underlying semantic patterns, leading to misclassification and unreliable predictions (Ye et al., 2025). This phenomenon shares a similar root cause with the amplification of hallucinations in generative tasks, but in the classification setting it manifests as unstable decision boundaries and degraded generalization (Almeida et al., 2021).

As a complementary paradigm, Reinforcement learning from human feedback (RLHF) (Stiennon et al., 2020; Ouyang et al., 2022) has demonstrated great potential in aligning LMs with human preferences, particularly in generative tasks such as dialogue and instruction following (Dong et al., 2024; Ouyang et al., 2025; Bai et al., 2022). Instead of predicting a single target, RLHF learns from preference signals and optimizes outputs through relative comparisons, thereby better reflecting subjective human judgments (Jiang et al., 2025). This naturally raises the question: can the principles of RLHF be transferred to classification tasks to overcome the inherent limitations of SFT?

Unfortunately, directly extending RLHF to classification faces fundamental challenges. First, unlike the rich sequential outputs of generative models, the output space of classification is discrete



(a) Implicit preference signals in human decisions.



(b) Improved text classification.

Figure 1: (a) An example illustrating the implicit preference signal in the construction of the SST-5 sentiment classification dataset. (b) The performance improvements of ClaHF on text classification tasks (e.g., higher accuracy and F1, lower error). More results are reported in Section 4.2.

and lacks expressive structure (Casper et al., 2023). Second, conventional RLHF relies heavily on manually annotated preference datasets, whose construction is costly and difficult to reuse across tasks and label spaces (Casper et al., 2023; Lindström et al., 2024). These issues make existing RLHF techniques hard to apply to classification, and relevant studies remain scarce.

To address these challenges, we propose ClaHF, a novel human feedback-inspired RL framework for text classification that introduces RLHF into classification training without requiring any additional human annotation. ClaHF is motivated by a key insight: although standard classification datasets do not contain explicit preference annotations, the ground-truth label of each instance implicitly encodes the preference signal underlying human decisions (Figure 1(a)). Based on this insight, ClaHF automatically constructs multiple candidate predictions and their ranking relations from raw classification data, yielding preference data at no extra cost. These preference signals are then used to train a RM that evaluates the relative quality of different predictions. Finally, the classification model is optimized with a preference-aligned proximal policy optimization (PPO) (Schulman et al., 2017) algorithm. Notably, ClaHF is applicable to both binary and multi-class classification.

Unlike conventional instance-wise supervision, ClaHF leverages a multi-stage collaborative process that enables the model not only to “predict correctly” but also to “align with preference signals” embedded in the data, thereby overcoming the limitations of SFT in classification. We conduct systematic experiments on eight benchmark datasets covering binary, multi-class, and code clas-

sification tasks. The results show that ClaHF significantly improves both performance and decision reliability over strong SFT baselines (Figure 1(b)). These findings demonstrate that ClaHF bridges the gap between RLHF and classification, and introduce a new perspective: classification models can benefit from preference alignment in addition to label supervision.

2 Related Work

Supervised Learning for Classification. The dominant paradigm is to perform SFT on labeled data using pre-trained models such as RoBERTa (Liu et al., 2019) and Qwen3 (Yang et al., 2025), which yields strong discriminative capability. Despite its empirical success, this training paradigm inherits the limitations of instance-wise supervision, which often leads to overconfidence, unreliable decisions, and misclassification of ambiguous samples near class boundaries. To alleviate these issues, a number of studies have proposed improved objectives. Lin et al. (2017) introduced focal loss to emphasize hard and minority samples, while Müller et al. (2019) proposed label smoothing to improve generalization by injecting soft targets and noise regularization. However, these methods still treat labels as independent supervision signals and fail to exploit the latent relative preferences and ranking relations among samples. Consequently, the training process remains essentially behavior cloning, which is insufficient to provide deeper guidance at the decision level.

Applications of RLHF. RLHF was originally developed to optimize agent behavior in simulated environments (Christiano et al., 2017), and was later applied to domains such as autonomous driving

(Wu et al., 2022) and robot learning (Wang et al., 2022). Prior work (Xu et al., 2023; Zhang et al., 2024; Fang et al., 2025) has also demonstrated the effectiveness of human feedback in image generation. More recently, the successful application of RLHF to large language models (LLMs) has substantially advanced NLP (Stiennon et al., 2020; Ouyang et al., 2022; Ziegler et al., 2019). By modeling human preferences and performing policy updates with PPO, RLHF integrates reward-driven mechanisms into language modeling, fundamentally reshaping dialogue reasoning, safety alignment, and long-form generation (Ji et al., 2025; Liu et al., 2023). However, within downstream NLP tasks, RLHF has so far been largely confined to generative scenarios, such as high-quality code test case generation (Steenhoek et al., 2025).

Preference Learning. Preference learning is a core component of RLHF and models relations among samples via pairwise or listwise comparisons. It has been widely adopted in ranking and recommendation systems (Fürnkranz and Hüllermeier, 2010). In NLP, preference learning is mainly used for generative tasks, although recent studies (Atapour-Abarghouei et al., 2021; Kim et al., 2023) have attempted to introduce preference modeling into classification to overcome the limitations of instance-wise supervision. For example, Kim et al. (2023) proposed the P2C method, which incorporates preference relations during training to enhance the model’s discrimination near class boundaries. Nevertheless, existing approaches typically treat preference ranking as an auxiliary loss jointly optimized with SFT; they neither construct an independent RM nor perform policy optimization with RL, and thus remain within the supervised learning paradigm.

3 Method

In this section, we present ClaHF, a novel human feedback-inspired framework for classification. Unlike conventional training paradigms that rely solely on instance-wise supervision, ClaHF borrows the core ideas of RLHF and introduces preference signals into text classification through reward modeling and RL, enabling finer calibration of the predictive distribution and improving both classification performance and decision reliability.

3.1 Problem Definition

We focus on text classification tasks, including both binary and multi-class settings. Given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, where x_i denotes an input sequence (typically a sentence or a paragraph) and y_i is the corresponding class label, conventional approaches learn a classifier f_θ by minimizing instance-wise losses to maximize the agreement between predictions and ground-truth labels. However, such instance-wise supervision provides only a single discriminative signal and cannot characterize the relative quality among different predictions, which limits further optimization.

To this end, we propose ClaHF, which augments standard label supervision with automatically constructed preference signals. Specifically, ClaHF introduces ranking constraints among multiple candidate predictions for the same instance to capture fine-grained differences in model outputs. The goal of ClaHF is therefore to learn a classifier f_θ that not only predicts the correct label but also optimizes its predictive distribution under preference constraints, thereby aligning the model with implicit “human feedback” embedded in the data. The overall framework of ClaHF is illustrated in Figure 2.

3.2 Supervised Fine-tuning

The SFT stage provides a stable and performant initial policy model f_θ for subsequent RL. In ClaHF, the model f_θ consists of a pre-trained LM g_ϕ and a randomly initialized classification head W_ψ . Given an input text x_i , we encode it into token sequences and feed them into f_θ to obtain hidden representations, followed by a softmax layer to produce the class probability distribution:

$$p_\theta(y_i | x_i) = \text{Softmax}(f_\theta(x_i)). \quad (1)$$

We optimize the model using the standard cross-entropy loss:

$$\mathcal{L}_{\text{sft}} = -\mathbb{E}_{(x_i, y_i) \sim \mathcal{D}} \log p_\theta(y_i | x_i). \quad (2)$$

This stage ensures that the policy model has strong baseline classification capability. However, the supervision signal is purely instance-wise and only enforces consistency with the ground-truth label, without reflecting the relative quality among alternative predictions.

3.3 Preference Data Construction

Training the RM requires data annotated with preference signals. In conventional RLHF, such data

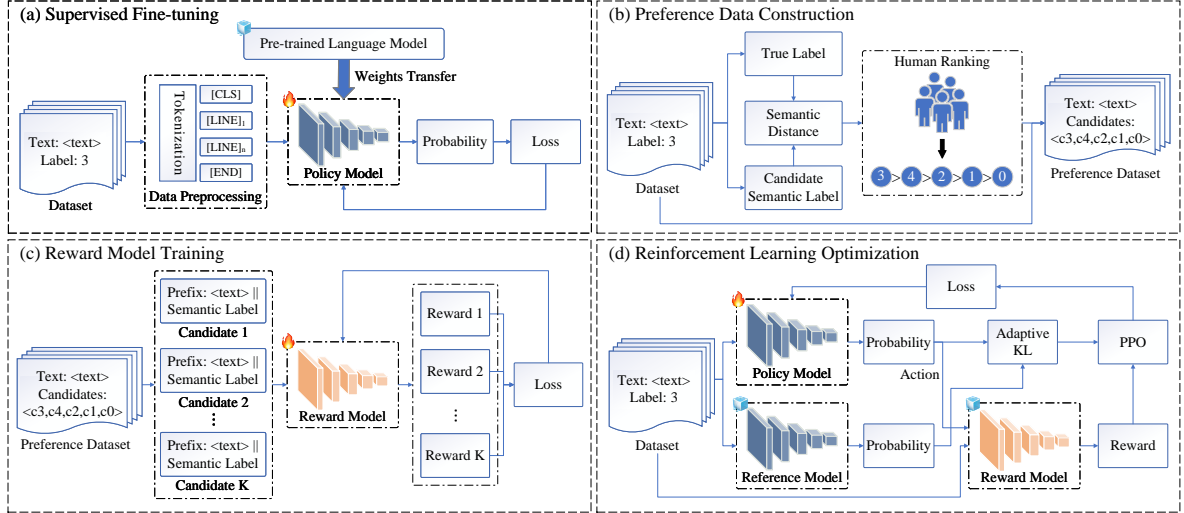


Figure 2: Overall framework of ClaHF. (a) SFT to provide high-quality initialization. (b) Automatic construction of preference data from the original classification dataset. (c) Training the RM with preference data. (d) RL optimization of the policy model using the trained RM.

Dataset	Label	Candidate Semantic Label
SST-5	0	Prediction: the sentence is very negative
	1	Prediction: the sentence is negative
	2	Prediction: the sentence is neutral
	3	Prediction: the sentence is positive
	4	Prediction: the sentence is very positive

Table 1: Ground-truth labels and corresponding candidate semantic labels for SST-5.

are manually labeled, which is expensive and task-specific. In contrast, classification datasets naturally contain ground-truth labels that implicitly encode human preferences. Based on this observation, ClaHF introduces a fully automated method for constructing preference data without any additional annotation.

Specifically, given a standard K -class dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, we first define a set of semantic label templates:

$$\mathcal{Y}_{\text{text}} = \{0 : c_0, 1 : c_1, \dots, K-1 : c_{K-1}\}, \quad (3)$$

where each c_j is a natural language description of class j , referred to as a candidate semantic label. Examples for SST-5 are shown in Table 1, and those for other datasets are provided in Appendix B.

For each instance (x_i, y_i) , we rank all candidate semantic labels c_j in ascending order of their semantic distance to the true label y_i , yielding a candidate list:

$$C_i = \{c_{i1}, c_{i2}, \dots, c_{iK}\}, \quad (4)$$

where c_{i1} always corresponds to the ground-truth class, and the remaining candidates are ordered by increasing semantic distance. For example, in SST-5, if the true label is 3, the resulting list is $\{c_3, c_4, c_2, c_1, c_0\}$.

Aggregating all instances produces the preference dataset $\mathcal{D}_{\text{pre}} = \{(x_i, C_i)\}_{i=1}^N$. Unlike the original dataset that only contains binary correctness information, \mathcal{D}_{pre} encodes relative plausibility among classes and naturally introduces a ‘‘Top-1 is better than others’’ preference structure.

3.4 Reward Model Training

The third stage of ClaHF trains a RM to evaluate the relative quality of text–candidate pairs and provide reward signals for subsequent policy optimization. The RM f_{rm} consists of the same pre-trained encoder g_ϕ as in SFT and a linear scoring head that outputs a scalar reward.

Given $\mathcal{D}_{\text{pre}} = \{(x_i, C_i)\}_{i=1}^N$, for each (x_i, C_i) we concatenate x_i with each c_{ij} , tokenize the sequence, and feed it into f_{rm} , which is expected to satisfy:

$$f_{\text{rm}}(x_i, c_{i1}) > \dots > f_{\text{rm}}(x_i, c_{iK}), \quad (5)$$

where c_{i1} corresponds to the ground-truth class.

Since classification lacks the rich candidate structure of generative tasks, we design a composite ranking loss with two components: $\mathcal{L}_{\text{top1}}$ and $\mathcal{L}_{\text{pairwise}}$. The Top-1 loss enforces the ground-truth

candidate to score higher than all others:

$$\mathcal{L}_{\text{top1}} = -\frac{1}{K-1} \sum_{j=2}^K \log \sigma(f_{\text{rm}}(x_i, c_{i1}) - f_{\text{rm}}(x_i, c_{ij})), \quad (6)$$

where σ is the sigmoid function.

The pairwise loss constrains the relative order among non-true candidates:

$$\mathcal{L}_{\text{pairwise}} = -\frac{2}{(K-1)(K-2)} \sum_{2 \leq j < k \leq K} \log \sigma(f_{\text{rm}}(x_i, c_{ij}) - f_{\text{rm}}(x_i, c_{ik})). \quad (7)$$

The overall objective is:

$$\mathcal{L}_{\text{rm}} = \alpha \mathcal{L}_{\text{top1}} + (1 - \alpha) \mathcal{L}_{\text{pairwise}}, \quad (8)$$

where $\alpha \in (0, 1]$ balances Top-1 preference and pairwise ranking consistency.

For binary classification, the loss simplifies to:

$$\mathcal{L}_{\text{rm}} = -\log \sigma(f_{\text{rm}}(x_i, c_{i1}) - f_{\text{rm}}(x_i, c_{i2})). \quad (9)$$

This training procedure enables the RM to capture finer-grained preference signals than conventional supervision.

3.5 Reinforcement Learning Optimization

In the final stage, we further optimize the policy model f_{θ} using PPO (Schulman et al., 2017) to improve decision quality and predictive reliability. A frozen reference model f_{ref} , with the same architecture as f_{θ} , is used to measure deviation from the supervised baseline.

Given an input x_i , the policy model f_{θ} outputs a class distribution $p_{\theta}(y_i | x_i)$, while the reference model f_{ref} outputs $p_{\text{ref}}(y_i | x_i)$. An action a_i is defined as a class label sampled from $p_{\theta}(y_i | x_i)$, and its raw reward $f_{\text{rm}}(x_i, a_i)$ is computed by the RM. To prevent the policy from drifting excessively from the supervised distribution, we introduce an adaptive KL-penalty (Ziegler et al., 2019) term into the reward and obtain the modified reward:

$$\hat{f}_{\text{rm}}(x_i, a_i) = f_{\text{rm}}(x_i, a_i) - \beta_i D_{\text{KL}}(p_{\theta}(\cdot | x_i) || p_{\text{ref}}(\cdot | x_i)), \quad (10)$$

where D_{KL} denotes the KL divergence between the two distributions. The KL coefficient β is dynamically adjusted by an adaptive controller:

$$\beta_{i+1} = \text{clip}(\beta_i \cdot (1 + \eta \cdot \frac{D_{\text{KL}} - D_{\text{target}}}{D_{\text{target}}}), \beta_{\text{min}}, \beta_{\text{max}}), \quad (11)$$

which keeps the KL divergence close to a target value D_{target} ; η is the update rate.

During PPO optimization, ClaHF maximizes the following clipped surrogate objective for each sample by comparing the probability ratio between the current and old policies:

$$\mathcal{L}_{\text{policy}} = \mathbb{E}_{x_i \sim \mathcal{D}}[-\min(r_i A_i, \text{clip}(r_i, 1 - \epsilon, 1 + \epsilon) A_i)], \quad (12)$$

where $r_i = \frac{p_{\theta}(a_i | x_i)}{p_{\text{ref}}(a_i | x_i)}$ is the probability ratio of the new and old policies, and A_i is the advantage function defined as:

$$A_i = \hat{f}_{\text{rm}}(x_i, a_i) - V(x_i), \quad (13)$$

with V denoting the output of the value function network. The value function is introduced to reduce training variance and stabilize policy learning. It is optimized using mean squared error:

$$\mathcal{L}_{\text{value}} = \mathbb{E}_{x_i \sim \mathcal{D}}[(V(x_i) - \hat{f}_{\text{rm}}(x_i, a_i))^2]. \quad (14)$$

The final optimization objective in the RL stage is:

$$\mathcal{L}_{\text{rl}} = \mathcal{L}_{\text{policy}} + \gamma \mathcal{L}_{\text{value}}, \quad (15)$$

where γ controls the relative weight of the value loss.

Through this design, ClaHF tightly couples preference-driven reward modeling with classification decisions, continuously refining decision boundaries while maintaining stability with respect to the policy model.

4 Experiments

4.1 Setups

Datasets. We select eight datasets to evaluate the effectiveness of ClaHF from three perspectives: binary classification, multi-class classification, and software engineering code tasks. Specifically, the binary classification tasks include: (1) SST-2 (Socher et al., 2013), (2) CoLA (Warstadt et al., 2019), and (3) MRPC (Dolan and Brockett, 2005). The multi-class classification tasks include: (4) AG News (Zhang et al., 2015), (5) SST-5 (Socher et al., 2013), and (6) Emotion (Saravia et al., 2018). The code tasks include: (7) Devign (Zhou et al., 2019), and (8) BigCloneBench (Wang et al., 2020). The data splits and additional details are provided in Appendix A.

Baselines. As existing RLHF frameworks are primarily designed for generative tasks and no mature solution is available for classification, we only

Backbone	Method	AG News			SST-5			Emotion		
		Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)	Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)	Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)
BERT	SFT	91.91 \pm 0.0	91.90 \pm 0.0	8.35 \pm 0.3	52.49 \pm 0.1	51.28 \pm 0.1	13.44 \pm 0.7	92.25 \pm 0.2	87.71 \pm 0.5	6.51 \pm 0.4
	ClaHF	92.18 \pm 0.1	92.17 \pm 0.1	6.96 \pm 0.5	54.30 \pm 0.1	51.50 \pm 0.8	11.36 \pm 0.3	92.85 \pm 0.1	88.05 \pm 0.4	5.74 \pm 0.5
RoBERTa	SFT	94.44 \pm 0.0	94.44 \pm 0.0	7.47 \pm 0.3	54.35 \pm 0.7	52.71 \pm 0.9	10.39 \pm 0.5	92.55 \pm 0.1	88.40 \pm 0.4	7.23 \pm 0.1
	ClaHF	94.61 \pm 0.0	94.60 \pm 0.0	5.52 \pm 0.2	55.87 \pm 0.4	54.17 \pm 0.6	8.43 \pm 0.3	92.93 \pm 0.1	88.93 \pm 0.1	6.80 \pm 0.1
T5	SFT	94.51 \pm 0.2	94.51 \pm 0.1	5.27 \pm 0.2	56.19 \pm 0.4	53.11 \pm 0.3	12.09 \pm 0.5	91.80 \pm 0.1	87.11 \pm 0.4	6.19 \pm 0.3
	ClaHF	94.79 \pm 0.1	94.79 \pm 0.1	4.52 \pm 0.1	57.47 \pm 0.8	54.72 \pm 0.9	11.02 \pm 0.4	92.20 \pm 0.1	88.11 \pm 0.1	4.74 \pm 0.3
OPT	SFT	93.77 \pm 0.1	93.77 \pm 0.1	9.60 \pm 0.5	52.85 \pm 0.4	49.72 \pm 0.5	10.42 \pm 0.8	92.28 \pm 0.4	87.74 \pm 0.7	6.62 \pm 0.1
	ClaHF	94.00 \pm 0.1	94.01 \pm 0.1	5.76 \pm 0.2	54.34 \pm 0.3	52.43 \pm 0.2	7.71 \pm 0.8	92.92 \pm 0.1	88.46 \pm 0.2	6.27 \pm 0.2
Qwen3	SFT	93.85 \pm 0.1	93.83 \pm 0.1	7.07 \pm 0.9	51.06 \pm 0.6	48.92 \pm 0.6	8.84 \pm 0.3	92.57 \pm 0.1	88.11 \pm 0.4	6.75 \pm 0.2
	ClaHF	94.17 \pm 0.1	94.16 \pm 0.1	5.22 \pm 0.5	52.67 \pm 0.1	49.55 \pm 1.0	7.61 \pm 0.4	93.62 \pm 0.1	90.18 \pm 0.1	5.06 \pm 0.2
Average		+0.25	+0.25	-1.96	+1.54	+1.32	-1.78	+0.61	+0.93	-0.94

Table 2: Test performance of SFT and ClaHF on three multi-class datasets. All values are reported as the mean and standard error over five random seeds.

compare ClaHF with standard SFT, where backbone models are trained on labeled data using cross-entropy loss. To verify the generality of ClaHF, we adopt nine widely used pre-trained LMs as backbones, including BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), OPT (Zhang et al., 2022), T5 (Raffel et al., 2020), and Qwen3 (Yang et al., 2025) for natural language tasks, as well as CodeBERT (Feng, 2020), CodeT5 (Wang et al., 2021), CodeT5+ (Wang et al., 2023), and CodeGen (Nijkamp et al., 2022) for code-related tasks.

Implementation details. All experiments are conducted on 8 NVIDIA RTX 4090 GPUs. Pre-trained weights are loaded via Hugging Face (Wolf et al., 2020). We use the AdamW optimizer (Loshchilov and Hutter, 2017) for all experiments and train each model for 10 epochs. During the SFT stage, the batch size is set to 16 and the learning rate to $2e-5$. For RM training, the learning rate is $1e-5$ and the preference loss weight α is set to 0.8. During the RL stage, the learning rate is $1e-6$ and the value loss coefficient γ is set to 0.25. Additional hyper-parameters and training details are reported in Appendix B.

4.2 Main Results

In this section, we systematically evaluate ClaHF on multi-class, binary, and code classification tasks and compare it with the corresponding SFT baselines across a wide range of backbone models, in order to assess its cross-model and cross-task generalization ability. To measure performance under class-imbalance scenarios, we report both Accuracy (Acc) and macro-F1 ($F1$). For CoLA, we report the commonly used Matthews Correlation Coefficient (MCC) (Chicco et al., 2021). In ad-

dition, we introduce Expected Calibration Error (ECE) (Guo et al., 2017) to evaluate predictive reliability.

4.2.1 Results on Multi-class Tasks

Table 2 summarizes the results of ClaHF and SFT on three multi-class datasets. ClaHF consistently outperforms SFT across all backbones and datasets, showing simultaneous improvements in Acc, F1, and ECE. This indicates that incorporating preference signals effectively complements instance-wise supervision and enables the model to better distinguish decision boundaries between different classes.

On the fine-grained sentiment dataset SST-5, where semantic differences are subtle, ClaHF achieves particularly notable gains, with an average improvement of 1.54 percentage points in Acc accompanied by a significant reduction in ECE. This demonstrates that ClaHF yields stronger discriminative ability and more reliable probability calibration in complex multi-class scenarios.

Similar trends are observed on the Emotion dataset. For example, with Qwen3, ClaHF improves Acc from 92.57% to 93.62% while reducing ECE from 6.75 to 5.06, indicating that ClaHF not only boosts accuracy but also fundamentally promotes more reliable and well-calibrated predictive distributions.

4.2.2 Results on Binary Tasks

Compared with multi-class classification, binary tasks are more sensitive to decision boundaries and impose stricter requirements on model confidence and distributional consistency. Therefore, we report these results separately.

Backbone	Method	SST-2			MRPC			CoLA	
		Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)	Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)	MCC (\uparrow)	ECE (\downarrow)
BERT	SFT	89.46 \pm 0.24	89.43 \pm 0.24	9.99 \pm 0.14	82.18 \pm 0.27	79.44 \pm 0.96	9.84 \pm 0.46	55.07 \pm 0.34	14.30 \pm 0.48
	ClaHF	90.10 \pm 0.09	90.09 \pm 0.10	9.04 \pm 0.55	83.28 \pm 0.29	80.53 \pm 0.29	8.64 \pm 0.21	56.10 \pm 0.24	12.16 \pm 0.57
RoBERTa	SFT	92.57 \pm 0.12	92.57 \pm 0.12	8.24 \pm 0.47	87.25 \pm 0.18	85.57 \pm 0.05	6.25 \pm 0.29	59.73 \pm 0.62	11.77 \pm 0.29
	ClaHF	93.19 \pm 0.23	93.19 \pm 0.23	6.40 \pm 0.13	87.86 \pm 0.15	86.28 \pm 0.23	5.12 \pm 0.20	61.21 \pm 0.82	8.43 \pm 0.37
T5	SFT	93.85 \pm 0.23	93.85 \pm 0.23	7.10 \pm 0.60	84.58 \pm 0.12	81.34 \pm 0.27	8.15 \pm 0.30	58.10 \pm 0.48	13.58 \pm 0.51
	ClaHF	94.63 \pm 0.04	94.64 \pm 0.04	5.03 \pm 0.13	85.78 \pm 0.28	83.68 \pm 0.28	6.83 \pm 0.27	59.23 \pm 0.46	11.75 \pm 0.33
OPT	SFT	90.52 \pm 0.46	90.51 \pm 0.45	8.76 \pm 0.52	83.61 \pm 0.25	80.69 \pm 0.37	7.55 \pm 0.19	56.42 \pm 0.26	13.95 \pm 0.77
	ClaHF	90.99 \pm 0.23	90.99 \pm 0.23	6.24 \pm 0.31	84.69 \pm 0.22	81.88 \pm 0.85	6.44 \pm 0.26	57.64 \pm 0.35	12.90 \pm 0.43
Qwen3	SFT	90.57 \pm 0.16	90.57 \pm 0.16	8.22 \pm 0.36	84.41 \pm 0.15	81.69 \pm 0.54	6.85 \pm 0.13	56.94 \pm 0.17	13.68 \pm 0.66
	ClaHF	91.47 \pm 0.29	91.47 \pm 0.29	7.09 \pm 0.30	85.02 \pm 0.18	82.97 \pm 0.26	5.58 \pm 0.36	58.30 \pm 0.27	12.44 \pm 0.34
Average		+0.68 \pm 0.07	+0.69 \pm 0.07	-1.70 \pm 0.29	+0.92 \pm 0.13	+1.32 \pm 0.27	-1.21 \pm 0.04	+1.24 \pm 0.08	-1.92 \pm 0.41

Table 3: Test performance of SFT and ClaHF on three binary datasets.

Backbone	Method	Devign			BigCloneBench		
		Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)	Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)
CodeBERT	SFT	62.15 \pm 0.19	58.77 \pm 0.50	13.48 \pm 0.81	97.12 \pm 0.14	97.11 \pm 0.14	4.01 \pm 0.43
	ClaHF	63.14 \pm 0.11	59.61 \pm 0.28	9.64 \pm 0.47	97.58 \pm 0.20	97.58 \pm 0.20	2.49 \pm 0.34
CodeT5	SFT	62.99 \pm 0.18	62.35 \pm 0.20	12.02 \pm 0.48	96.80 \pm 0.17	96.79 \pm 0.17	4.93 \pm 0.83
	ClaHF	63.49 \pm 0.36	62.50 \pm 0.56	9.06 \pm 0.24	97.29 \pm 0.16	97.29 \pm 0.16	2.46 \pm 0.26
CodeT5+	SFT	63.20 \pm 0.17	62.87 \pm 0.25	11.13 \pm 0.42	97.87 \pm 0.10	97.87 \pm 0.10	2.08 \pm 0.06
	ClaHF	64.04 \pm 0.15	63.57 \pm 0.14	8.60 \pm 0.53	98.14 \pm 0.07	98.13 \pm 0.07	1.79 \pm 0.10
CodeGen	SFT	62.38 \pm 0.29	59.29 \pm 0.24	15.11 \pm 0.44	97.60 \pm 0.14	97.60 \pm 0.14	2.43 \pm 0.27
	ClaHF	63.45 \pm 0.11	61.95 \pm 0.87	11.13 \pm 0.71	97.98 \pm 0.09	97.98 \pm 0.09	1.92 \pm 0.05
Average		+0.85 \pm 0.12	+1.09 \pm 0.55	-3.33 \pm 0.35	+0.40 \pm 0.05	+0.40 \pm 0.05	-1.20 \pm 0.50

Table 4: Test performance of SFT and ClaHF on code datasets.

As shown in Table 3, ClaHF consistently outperforms SFT on SST-2, MRPC, and CoLA across all backbones, with particularly large gains in ECE. On SST-2, ClaHF achieves simultaneous improvements in Acc and F1, demonstrating that preference signals help the model better capture the relative ordering between positive and negative samples.

The advantage of ClaHF is even more pronounced on MRPC, which requires fine-grained judgments of semantic equivalence and is highly sensitive to boundary cases. Compared with SFT, ClaHF substantially improves performance and reduces ECE by about 1.21, indicating enhanced decision stability near hard examples. On CoLA, ClaHF improves MCC by approximately 1.24 while reducing ECE by about 1.92, further confirming its dual benefits in refining decision boundaries and improving confidence calibration.

4.2.3 Results on Code Tasks

Compared with natural language, source code is a structured and formal language with strict syntactic and semantic constraints, which places higher demands on model discrimination and stability.

Table 4 presents the results on Devign and

BigCloneBench. Overall, ClaHF achieves consistent improvements across all code backbones. On Devign, ClaHF reduces ECE by an average of 3.33 while also yielding small but consistent gains in Acc and F1, demonstrating its effectiveness on structured code data. On BigCloneBench, ClaHF again shows stable advantages, especially in terms of ECE, indicating that it maintains reliable confidence calibration even in high-accuracy regimes—an essential property for practical software engineering applications.

4.3 Ablation Study

To further investigate the contribution of each stage in ClaHF, we conduct ablation studies from two perspectives. Unless otherwise specified, all experiments in this section are conducted on the Emotion dataset (Saravia et al., 2018) with the Qwen3 (Yang et al., 2025) backbone, and all other settings are kept identical.

4.3.1 Effect of RM Training Objectives

Table 5 reports the impact of different RM training objectives on the final classification performance. In this group of experiments, we only vary the

Method	\mathcal{L}_{sft}	$\mathcal{L}_{\text{top1}}$	$\mathcal{L}_{\text{pairwise}}$	Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)
SFT	✓	–	–	92.57 \pm 0.14	88.11 \pm 0.40	6.75 \pm 0.22
ClaHF-Pair	✓	–	✓	92.78 \pm 0.14	88.88 \pm 0.12	6.44 \pm 0.21
ClaHF-Top1	✓	✓	–	93.15 \pm 0.20	89.28 \pm 0.40	5.98 \pm 0.15
ClaHF (Full)	✓	✓	✓	93.62 \pm 0.10	90.18 \pm 0.12	5.06 \pm 0.21

Table 5: Effect of RM training objectives on classification performance.

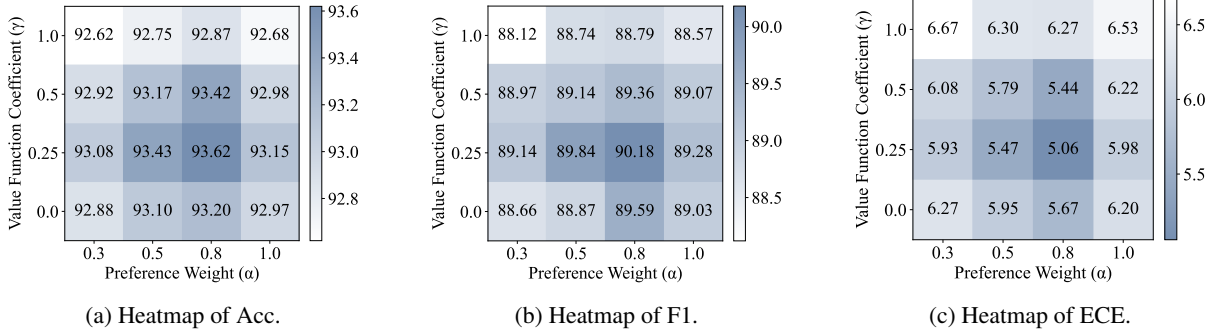


Figure 3: Performance heatmaps of ClaHF in terms of Acc, F1, and ECE over a continuous $\alpha \times \gamma$ grid.

loss components used in the RM stage. To isolate the effect of preference signals, we construct two variants of ClaHF: (1) ClaHF-Pair: applies equal constraints to all candidates. (2) ClaHF-Top1: only enforces the margin between the Top-1 candidate and the remaining candidates.

As shown in the table, all three ClaHF variants outperform SFT, confirming that preference signals effectively complement single-label supervision. However, ClaHF-Pair yields only limited improvements, indicating that enforcing pairwise ranking consistency alone is insufficient to provide strong discriminative signals and behaves similarly to conventional learning-to-rank methods.

When both Top-1 and pairwise constraints are combined, the model achieves the best performance across Acc, F1, and ECE, which significantly exceeds SFT and the two variants. This result indicates that the two types of preference signals are complementary: the Top-1 constraint emphasizes discriminability of the optimal decision, while the pairwise constraint preserves the global ranking structure. Additional studies are reported in Appendix C.

4.3.2 Stability Analysis

We further analyze the joint effect of the preference weight α and the value coefficient γ over a continuous hyperparameter space to assess the stability of ClaHF and the structure of its optimal region.

Figure 3 presents heatmaps of Acc, F1, and ECE over a continuous $\alpha \times \gamma$ grid. The performance sur-

faces vary smoothly rather than exhibiting abrupt oscillations, indicating that ClaHF is robust to hyperparameter perturbations and that small changes within reasonable ranges do not lead to dramatic performance drops.

Moreover, the optimal performance does not concentrate at a single point but spans a coherent region. Specifically, when $\alpha \in [0.7, 0.9]$ and $\gamma \in [0.2, 0.3]$, the model achieves near-optimal Acc and F1 while maintaining low ECE. The existence of this structured optimal region demonstrates that the improvements brought by ClaHF are not the result of accidental tuning, but remain stable across a range of reasonable configurations.

5 Conclusion

We propose ClaHF, a human-feedback-inspired RL framework for text classification, which for the first time systematically introduces the complete RLHF pipeline into classification tasks. ClaHF is tailored to the decision structure of classification and incorporates dedicated designs for preference data construction, reward modeling, and policy optimization, enabling RL to effectively optimize discriminative prediction objectives. Compared with SFT, ClaHF achieves stable and consistent improvements in both performance and calibration across a wide range of natural language and code classification tasks, demonstrating that preference-driven RL is also a powerful paradigm for enhancing the effectiveness and reliability of classification models.

557 Limitations

558 Despite the promising results of ClaHF on multi-
559 ple text classification benchmarks, it has two main
560 limitations. (1) ClaHF relies on automatically con-
561 structed preference data derived from original la-
562 bels. Although this avoids additional human anno-
563 tation costs, the quality of RM training is inevitably
564 constrained by the noise and bias in the original
565 datasets. Integrating ClaHF with real human feed-
566 back is a promising direction to further improve
567 performance and robustness. (2) Our experiments
568 focus on single-label text classification. The appli-
569 cability of ClaHF to more complex discriminative
570 settings, such as multi-label and hierarchical classi-
571 fication, remains unexplored and is left for future
572 work.

573 Acknowledgments

574 References

575 Matthew Almeida, Yong Zhuang, Wei Ding, Scott E
576 Crouter, and Ping Chen. 2021. Mitigating class-
577 boundary label uncertainty to reduce both model bias
578 and variance. *ACM Transactions on Knowledge Dis-*
579 *covery from Data (TKDD)*, 15(2):1–18.

580 Amir Atapour-Abarghouei, Stephen Bonner, and An-
581 drew Stephen McGough. 2021. Rank over class: The
582 untapped potential of ranking in natural language pro-
583 cessing. In *2021 IEEE International Conference on*
584 *Big Data (Big Data)*, pages 3950–3959. IEEE.

585 Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda
586 Askell, Anna Chen, Nova DasSarma, Dawn Drain,
587 Stanislav Fort, Deep Ganguli, Tom Henighan, and 1
588 others. 2022. Training a helpful and harmless assis-
589 tant with reinforcement learning from human feed-
590 back. *arXiv preprint arXiv:2204.05862*.

591 Stephen Casper, Xander Davies, Claudia Shi,
592 Thomas Krendl Gilbert, Jérémy Scheurer, Javier
593 Rando, Rachel Freedman, Tomasz Korbak, David
594 Lindner, Pedro Freire, and 1 others. 2023. Open
595 problems and fundamental limitations of reinforcement
596 learning from human feedback. *arXiv preprint*
597 *arXiv:2307.15217*.

598 Davide Chicco, Matthijs J Warrens, and Giuseppe Ju-
599 rman. 2021. The matthews correlation coefficient
600 (mcc) is more informative than cohen’s kappa and
601 brier score in binary classification assessment. *Ieee*
602 *Access*, 9:78368–78381.

603 Paul F Christiano, Jan Leike, Tom Brown, Miljan Mar-
604 tic, Shane Legg, and Dario Amodei. 2017. Deep
605 reinforcement learning from human preferences. *Ad-*
606 *vances in neural information processing systems*, 30.

607 Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Sheng-
608 bang Tong, Saining Xie, Dale Schuurmans, Quoc V

Le, Sergey Levine, and Yi Ma. 2025. Sft mem- 609
orizes, rl generalizes: A comparative study of 610
foundation model post-training. *arXiv preprint* 611
arXiv:2501.17161. 612

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 613
Kristina Toutanova. 2019. Bert: Pre-training of deep 614
bidirectional transformers for language understand- 615
ing. In *Proceedings of the 2019 conference of the* 616
North American chapter of the association for com- 617
putational linguistics: human language technologies, 618
volume 1 (long and short papers), pages 4171–4186. 619

William B Dolan and Chris Brockett. 2005. Automati- 620
cally constructing a corpus of sentential paraphrases. 621
In *Proceedings of the third international workshop* 622
on paraphrasing (IWP2005). 623

Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, 624
Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, 625
Caiming Xiong, and Tong Zhang. 2024. Rlhf work- 626
flow: From reward modeling to online rlhf. *arXiv* 627
preprint arXiv:2405.07863. 628

Tiyu Fang, Mingxin Zhang, Ran Song, Xiaolei Li, 629
Zhiyuan Wei, and Wei Zhang. 2025. Human-guided 630
data augmentation via diffusion model for surface 631
defect recognition under limited data. *IEEE Transac-* 632
tions on Instrumentation and Measurement. 633

Z Feng. 2020. Codebert: A pre-trained model for 634
program-ming and natural languages. *arXiv preprint* 635
arXiv:2002.08155. 636

Guilherme Fonseca, Washington Cunha, Gabriel Pre- 637
nassi, Marcos André Gonçalves, and Leonardo 638
Chaves Dutra Da Rocha. 2025. Instance-selection- 639
inspired undersampling strategies for bias reduction 640
in small and large language models for binary text 641
classification. In *Proceedings of the 63rd Annual* 642
Meeting of the Association for Computational Lin- 643
guistics (Volume 1: Long Papers), pages 9323–9340. 644

Zihao Fu, Anthony Man-Cho So, and Nigel Collier. 645
2023. A stability analysis of fine-tuning a pre-trained 646
model. *arXiv preprint arXiv:2301.09820*. 647

Johannes Fürnkranz and Eyke Hüllermeier. 2010. Pref- 648
erence learning and ranking by pairwise comparison. 649
In *Preference learning*, pages 65–82. Springer. 650

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Wein- 651
berger. 2017. On calibration of modern neural net- 652
works. In *International conference on machine learn-* 653
ing, pages 1321–1330. PMLR. 654

Miaomiao Ji, Yanqiu Wu, Zhibin Wu, Shoujin Wang, 655
Jian Yang, Mark Dras, and Usman Naseem. 2025. 656
A survey on progress in llm alignment from the 657
perspective of reward design. *arXiv preprint* 658
arXiv:2505.02666. 659

Ruili Jiang, Kehai Chen, Xuefeng Bai, Zhixuan He, 660
Juntao Li, Muyun Yang, Tiejun Zhao, Liqiang Nie, 661
and Min Zhang. 2025. A survey on human preference 662
learning for aligning large language models. *ACM* 663
Computing Surveys, 58(6). 664

665	Jaehyung Kim, Jinwoo Shin, and Dongyeop Kang. 2023.	Sheng Ouyang, Yulan Hu, Ge Chen, Qingyang Li,	718
666	Prefer to classify: Improving text classifiers via auxil-	Fuzheng Zhang, and Yong Liu. 2025. Towards re-	719
667	ary preference learning. In <i>International Conference</i>	ward fairness in rlhf: From a resource allocation	720
668	on <i>Machine Learning</i> , pages 16807–16828. PMLR.	perspective. <i>arXiv preprint arXiv:2505.23349</i> .	721
669	Qian Li, Hao Peng, Jianxin Li, Congying Xia, Renyu	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	722
670	Yang, Lichao Sun, Philip S Yu, and Lifang He. 2022.	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	723
671	A survey on text classification: From traditional to	Wei Li, and Peter J Liu. 2020. Exploring the lim-	724
672	deep learning. <i>ACM Transactions on Intelligent Sys-</i>	its of transfer learning with a unified text-to-text	725
673	<i>tems and Technology (TIST)</i> , 13(2):1–41.	transformer. <i>Journal of machine learning research</i> ,	726
674	Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He,	21(140):1–67.	727
675	and Piotr Dollár. 2017. Focal loss for dense object	Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang,	728
676	detection. In <i>Proceedings of the IEEE international</i>	Junlin Wu, and Yi-Shin Chen. 2018. Carer: Con-	729
677	<i>conference on computer vision</i> , pages 2980–2988.	textualized affect representations for emotion recog-	730
678	Adam Dahlgren Lindström, Leila Methnani, Lea	nition. In <i>Proceedings of the 2018 conference on</i>	731
679	Krause, Petter Ericson, Íñigo Martínez de Rituerto	<i>empirical methods in natural language processing</i> ,	732
680	de Troya, Dimitri Coelho Mollo, and Roel Dobbe.	pages 3687–3697.	733
681	2024. Ai alignment through reinforcement learning	John Schulman, Filip Wolski, Prafulla Dhariwal,	734
682	from human feedback? contradictions and limita-	Alec Radford, and Oleg Klimov. 2017. Proxi-	735
683	tions. <i>arXiv preprint arXiv:2406.18346</i> .	mal policy optimization algorithms. <i>arXiv preprint</i>	736
684	Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang,	<i>arXiv:1707.06347</i> .	737
685	Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li,	Richard Socher, Alex Perelygin, Jean Wu, Jason	738
686	Mengshen He, Zhengliang Liu, and 1 others. 2023.	Chuang, Christopher D Manning, Andrew Y Ng, and	739
687	Summary of chatgpt-related research and perspective	Christopher Potts. 2013. Recursive deep models for	740
688	towards the future of large language models. <i>Meta-</i>	semantic compositionality over a sentiment treebank.	741
689	<i>radiology</i> , 1(2):100017.	In <i>Proceedings of the 2013 conference on empiri-</i>	742
690	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	<i>cal methods in natural language processing</i> , pages	743
691	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	1631–1642.	744
692	Luke Zettlemoyer, and Veselin Stoyanov. 2019.	Benjamin Steenhoeck, Michele Tufano, Neel Sundare-	745
693	Roberta: A robustly optimized bert pretraining ap-	san, and Alexey Svyatkovskiy. 2025. Reinforcement	746
694	proach. <i>arXiv preprint arXiv:1907.11692</i> .	learning from automatic feedback for high-quality	747
695	Ilya Loshchilov and Frank Hutter. 2017. Decou-	unit test generation. In <i>2025 IEEE/ACM Interna-</i>	748
696	pled weight decay regularization. <i>arXiv preprint</i>	<i>tional Workshop on Deep Learning for Testing and</i>	749
697	<i>arXiv:1711.05101</i> .	<i>Testing for Deep Learning (DeepTest)</i> , pages 37–44.	750
698	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey	IEEE.	751
699	Svyatkovskiy, Ambrosio Blanco, Colin Clement,	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel	752
700	Dawn Drain, Daxin Jiang, Duyu Tang, and 1 others.	Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,	753
701	2021. Codexglue: A machine learning benchmark	Dario Amodei, and Paul F Christiano. 2020. Learn-	754
702	dataset for code understanding and generation. <i>arXiv</i>	ing to summarize with human feedback. <i>Advances</i>	755
703	<i>preprint arXiv:2102.04664</i> .	<i>in neural information processing systems</i> , 33:3008–	756
704	Rafael Müller, Simon Kornblith, and Geoffrey E Hinton.	3021.	757
705	2019. When does label smoothing help? <i>Advances</i>	Xiaofei Sun, Xiaoya Li, Jiwei Li, Fei Wu, Shang-	758
706	<i>in neural information processing systems</i> , 32.	wei Guo, Tianwei Zhang, and Guoyin Wang. 2023.	759
707	Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan	Text classification via large language models. <i>arXiv</i>	760
708	Wang, Yingbo Zhou, Silvio Savarese, and Caiming	<i>preprint arXiv:2305.08377</i> .	761
709	Xiong. 2022. Codegen: An open large language	Alex Wang, Amanpreet Singh, Julian Michael, Felix	762
710	model for code with multi-turn program synthesis.	Hill, Omer Levy, and Samuel Bowman. 2018. Glue:	763
711	<i>arXiv preprint arXiv:2203.13474</i> .	A multi-task benchmark and analysis platform for	764
712	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	natural language understanding. In <i>Proceedings of</i>	765
713	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	<i>the 2018 EMNLP workshop BlackboxNLP: Analyzing</i>	766
714	Sandhini Agarwal, Katarina Slama, Alex Ray, and 1	<i>and interpreting neural networks for NLP</i> , pages 353–	767
715	others. 2022. Training language models to follow in-	355.	768
716	structions with human feedback. <i>Advances in neural</i>	Wenhan Wang, Ge Li, Bo Ma, Xin Xia, and Zhi	769
717	<i>information processing systems</i> , 35:27730–27744.	Jin. 2020. Detecting code clones with graph neu-	770
		ral network and flow-augmented abstract syntax	771
		tree. In <i>2020 IEEE 27th International Conference</i>	772
		<i>on Software Analysis, Evolution and Reengineering</i>	773
		<i>(SANER)</i> , pages 261–271. IEEE.	774

775	Xiaofei Wang, Kimin Lee, Kourosh Hakhmaneshi, Pieter Abbeel, and Michael Laskin. 2022. Skill preferences: Learning to extract and execute robotic skills from human feedback. In <i>Conference on robot learning</i> , pages 1259–1268. PMLR.	832
776		833
777		834
778		835
779		836
780	Yue Wang, Hung Le, Akhilesh Gotmare, Nghi Bui, Junnan Li, and Steven Hoi. 2023. Codet5+: Open code large language models for code understanding and generation. In <i>Proceedings of the 2023 conference on empirical methods in natural language processing</i> , pages 1069–1088.	837
781		838
782		
783		
784		
785		
786	Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. 2021. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. <i>arXiv preprint arXiv:2109.00859</i> .	
787		
788		
789		
790		
791	Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. 2019. Neural network acceptability judgments. <i>Transactions of the Association for Computational Linguistics</i> , 7:625–641.	
792		
793		
794		
795	Bowen Wei and Ziwei Zhu. 2025. Protolens: Advancing prototype learning for fine-grained interpretability in text classification. In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4503–4523.	
796		
797		
798		
799		
800		
801	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, and 1 others. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations</i> , pages 38–45.	
802		
803		
804		
805		
806		
807		
808		
809	Jingda Wu, Zhiyu Huang, Wenhui Huang, and Chen Lv. 2022. Prioritized experience-based reinforcement learning with human guidance for autonomous driving. <i>IEEE Transactions on Neural Networks and Learning Systems</i> , 35(1):855–869.	
810		
811		
812		
813		
814	Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. 2023. Imagereward: Learning and evaluating human preferences for text-to-image generation. <i>Advances in Neural Information Processing Systems</i> , 36:15903–15935.	
815		
816		
817		
818		
819		
820	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	
821		
822		
823		
824		
825	Junjie Ye, Yuming Yang, Yang Nan, Shuo Li, Qi Zhang, Tao Gui, Xuan-Jing Huang, Peng Wang, Zhongchao Shi, and Jianping Fan. 2025. Analyzing the effects of supervised fine-tuning on model knowledge from token and parameter levels. In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 471–513.	
826		
827		
828		
829		
830		
831		
	Shu Zhang, Xinyi Yang, Yihao Feng, Can Qin, Chia-Chih Chen, Ning Yu, Zeyuan Chen, Huan Wang, Silvio Savarese, Stefano Ermon, and 1 others. 2024. Hive: Harnessing human feedback for instructional visual editing. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 9026–9036.	839
		840
		841
		842
		843
	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, and 1 others. 2022. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> .	
	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. <i>Advances in neural information processing systems</i> , 28.	844
		845
		846
		847
	Yaqin Zhou, Shangqing Liu, Jingkai Siow, Xiaoning Du, and Yang Liu. 2019. Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks. <i>Advances in neural information processing systems</i> , 32.	848
		849
		850
		851
		852
	Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. <i>arXiv preprint arXiv:1909.08593</i> .	853
		854
		855
		856
		857
	A Datasets	858
	As described in Section 4.1, we evaluate ClaHF on eight representative text classification benchmarks. For all datasets, we set the maximum input length to 400 tokens for tokenization. The detailed descriptions are as follows.	859
		860
		861
		862
		863
	SST-2 (Socher et al., 2013) is a single-sentence binary classification dataset consisting of sentences from movie reviews annotated with sentiment labels. The labels are 0 and 1, where 0 denotes negative and 1 denotes positive. The dataset contains 67,350 training samples, 873 validation samples, and 1,821 test samples. As part of the GLUE benchmark (Wang et al., 2018), the official dataset is available at https://huggingface.co/datasets/nyu-ml1/glue/viewer/sst2 .	864
		865
		866
		867
		868
		869
		870
		871
		872
		873
	CoLA (Warstadt et al., 2019) is a single-sentence binary classification dataset for linguistic acceptability. The corpus is collected from books and journals in linguistic theory. Each sentence is annotated as to whether it is grammatically acceptable, with labels 0 and 1, where 0 indicates ungrammatical and 1 indicates grammatical. The dataset contains 8,551 training samples, 1,043 validation samples, and 1,063 test samples. It is also part of the GLUE benchmark (Wang et al.,	874
		875
		876
		877
		878
		879
		880
		881
		882
		883

884	2018) and is available at https://huggingface.co/datasets/nyu-ml/glue/viewer/cola .	936
885		937
886	MRPC (Dolan and Brockett, 2005) is a binary	938
887	classification dataset consisting of sentence pairs	939
888	extracted from online news sources, manually an-	940
889	notated to indicate whether the two sentences are	941
890	semantically equivalent. The labels are 0 and 1,	942
891	where 0 denotes not paraphrase and 1 denotes	943
892	paraphrase. The dataset contains 3,668 training	944
893	samples, 408 validation samples, and 1,725 test	945
894	samples. It is also part of GLUE (Wang et al.,	946
895	2018) and is available at https://huggingface.co/datasets/nyu-ml/glue/viewer/mrpc .	947
896		948
897	AG News (Zhang et al., 2015) is a single-	949
898	sentence four-class news topic classification dataset	
899	collected from global news articles. The dataset	
900	covers four major categories: World (label 0),	
901	Sports (label 1), Business (label 2), and Sci-	
902	ence/Technology (label 3). It contains 120,000	
903	training samples and 7,600 test samples. Since no	
904	validation split is provided, we randomly sample	
905	20% of the training set as the validation set. The	
906	dataset is available at https://huggingface.co/datasets/fancyzhx/ag_news .	
907		
908	SST-5 (Socher et al., 2013) is a fine-grained	
909	single-sentence five-class sentiment classification	
910	dataset for movie reviews. The labels range	
911	from 0 to 4, representing very negative, nega-	
912	tive, neutral, positive, and very positive, respec-	
913	tively. The dataset contains 8,544 training sam-	
914	ples, 1,101 validation samples, and 2,210 test sam-	
915	ples. It is available at https://huggingface.co/datasets/SetFit/sst5 .	
916		
917	Emotion (Saravia et al., 2018) is a single-	
918	sentence six-class emotion classification dataset	
919	built from English Twitter messages. It covers six	
920	emotion categories: sadness (label 0), joy (label 1),	
921	love (label 2), anger (label 3), fear (label 4), and sur-	
922	prise (label 5). The dataset contains 16,000 training	
923	samples, 2,000 validation samples, and 2,000 test	
924	samples. It is available at https://huggingface.co/datasets/dair-ai/emotion .	
925		
926	Devign (Zhou et al., 2019) is a code defect	
927	detection dataset designed to identify whether	
928	a code snippet is potentially vulnerable to soft-	
929	ware attacks. The labels are 0 and 1, where	
930	0 denotes safe code and 1 denotes vulnerable	
931	code. Following the CodeXGLUE (Lu et al.,	
932	2021), we split the dataset into training, valida-	
933	tion, and test sets with an 8:1:1 ratio, resulting	
934	in 21,854 training samples and 2,732 samples for	
935	both validation and test sets. The dataset is avail-	
	able at https://huggingface.co/datasets/google/code_x_glue_cc_defect_detection .	936
		937
	BigCloneBench (Wang et al., 2020) is a	938
	manually annotated Java code clone detection	939
	benchmark for code clone detection and semantic	940
	similarity learning. The labels are 0 and 1, where	941
	1 indicates that the code pair is semantically	942
	equivalent and 0 indicates non-equivalence.	943
	Since the original dataset is very large, we	944
	follow the CodeXGLUE (Lu et al., 2021) setting	945
	and use only 10% of the data for training and	946
	validation. The dataset is available at https://huggingface.co/datasets/google/code_x_glue_cc_clone_detection_big_clone_bench .	947
		948
		949
	B Additional Implementation Details	950
		951
	In the preference data construction stage, for each	952
	ground-truth class in the eight datasets, we define	953
	a natural-language candidate semantic label, as	954
	shown in Table 6. The core motivation is to map	955
	the original discrete class label y_i into a semantic	956
	description c_j that can be directly interpreted by	957
	LMs.	958
	In addition, during RM training, we concatenate	959
	each input sample with its candidate semantic label.	960
	Different datasets use different prefixes to explicitly	961
	distinguish task contexts and align the inputs with	962
	their semantic labels. The unified input format	963
	for the RM is $Prefix : x_i c_{ij}$, the prefixes for	964
	different datasets are defined as follows:	965
		966
	• AG News: News	967
		968
	• SST-5 / SST-2 / CoLA / MRPC: Sentence	969
		970
	• Emotion: Text	971
		972
	• Devign / BigCloneBench: Code	973
		974
		975
		976
	C Additional Ablation Studies	977
		978
	This section further investigates ablation studies	979
	of ClaHF. Following Section 4.3, all experiments	980
	are conducted on the Emotion dataset with Qwen3	981
	as the backbone model. All reported values are	982
	the mean and standard error over five runs with	
	different random seeds.	
	Effect of the Value Function Loss in RL. Ta-	
	ble 7 analyzes the influence of the value loss weight	
	in the RL stage. In these experiments, all RM set-	
	tings are fixed and α is set to 0.8, while only the	
	value loss coefficient γ is varied.	
	Without the value loss (ClaHF-NoValue), the	
	model still outperforms SFT, showing that	

Dataset \mathcal{D}	True Label y_i	Candidate Semantic Label c_j
AG News	0	Prediction: the news topic is World
	1	Prediction: the news topic is Sports
	2	Prediction: the news topic is Business
	3	Prediction: the news topic is Science/Technology
SST-5	0	Prediction: the sentence is very negative
	1	Prediction: the sentence is negative
	2	Prediction: the sentence is neutral
	3	Prediction: the sentence is positive
	4	Prediction: the sentence is very positive
Emotion	0	Prediction: the emotion in the text is sadness
	1	Prediction: the emotion in the text is joy
	2	Prediction: the emotion in the text is love
	3	Prediction: the emotion in the text is anger
	4	Prediction: the emotion in the text is fear
	5	Prediction: the emotion in the text is surprise
SST-2	0	Prediction: the sentence is negative
	1	Prediction: the sentence is positive
CoLA	0	Prediction: the sentence is grammatically unacceptable
	1	Prediction: the sentence is grammatically acceptable
MRPC	0	Prediction: the two sentences are not semantically equivalent
	1	Prediction: the two sentences are semantically equivalent
Devign	0	Prediction: the code is non-vulnerable
	1	Prediction: the code is vulnerable
MRPC	0	Prediction: the two code snippets are not semantically similar
	1	Prediction: the two code snippets are semantically similar

Table 6: Ground-truth labels and their corresponding candidate semantic labels for the eight datasets.

Method	$\mathcal{L}_{\text{value}}$	Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)
NoValue	–	93.20 \pm 0.15	89.59 \pm 0.45	5.67 \pm 0.12
ClaHF	$\gamma = 0.25$	93.62 \pm 0.10	90.18 \pm 0.12	5.06 \pm 0.21
ClaHF	$\gamma = 0.5$	93.42 \pm 0.13	89.36 \pm 0.26	5.44 \pm 0.12
ClaHF	$\gamma = 1.0$	92.87 \pm 0.17	88.79 \pm 0.34	6.27 \pm 0.34

Table 7: Effect of the value loss weight during RL.

983 preference-driven policy optimization is effective
984 by itself. However, its performance is clearly infe-
985 rior to the full ClaHF, indicating that RL without
986 value function regularization is more susceptible to
987 high-variance reward signals.

988 When a moderate value loss is introduced, the
989 model achieves the best results on Acc, F1, and
990 ECE, validating the crucial role of the value func-
991 tion in stabilizing training and smoothing policy
992 updates. As γ further increases, performance starts
993 to degrade, suggesting that overly strong value con-
994 straints hinder effective preference-driven updates
995 and thus weaken the benefits of ClaHF.

996 Table 8 reports a comparison between ClaHF-
997 Pair with different value function loss weights γ
998 and the SFT baseline. In this experiment, we keep
999 the RM and all other PPO hyperparameters fixed
1000 and only vary the weight of the value function loss
1001 to analyze the model behavior in the absence of
1002 the Top-1 constraint, i.e., when equal constraints
1003 are imposed on all candidates without explicitly
1004 emphasizing the optimal candidate and the model
1005 is trained only to learn the relative ordering among
1006 candidates.

1007 From the results, we observe that when the value
1008 function loss is not introduced and without the Top-
1009 1 constraint, the model can hardly outperform the
1010 SFT baseline on Acc, F1, and ECE, and even ex-
1011 hibits slight degradation on the calibration metric
1012 ECE. This indicates that relying solely on pair-
1013 wise preference signals for policy updates makes
1014 it difficult for the model to stably learn effective
1015 discriminative boundaries.

1016 After introducing a moderate value function loss,
1017 ClaHF-Pair achieves modest improvements over
1018 SFT, which further verifies that the value function
1019 can alleviate the instability caused by high-variance
1020 rewards in preference-based policy optimization.
1021 However, as γ increases further, the model perfor-
1022 mance drops and even falls below that of SFT, sug-
1023 gesting that an excessively large value loss weight
1024 overly suppresses policy updates when the Top-1
1025 constraint is absent.

Method	$\mathcal{L}_{\text{value}}$	Acc (\uparrow)	F1 (\uparrow)	ECE (\downarrow)
SFT	–	92.57 \pm 0.14	88.11 \pm 0.40	6.75 \pm 0.22
ClaHF-Pair	–	92.53 \pm 0.13	88.26 \pm 0.23	6.92 \pm 0.38
ClaHF-Pair	$\gamma = 0.25$	92.78 \pm 0.14	88.88 \pm 0.12	6.44 \pm 0.21
ClaHF-Pair	$\gamma = 0.5$	92.60 \pm 0.10	88.57 \pm 0.09	6.69 \pm 0.28
ClaHF-Pair	$\gamma = 1.0$	92.28 \pm 0.14	87.27 \pm 0.18	6.97 \pm 0.22

Table 8: Effect of the value loss weight on ClaHF-Pair.

1026 In summary, these results demonstrate that in
1027 ClaHF-Pair, although the value function loss can
1028 improve training stability, its effective range is
1029 clearly limited and requires careful tuning. This
1030 finding also provides additional evidence for our
1031 main conclusion that Top-1 and pairwise prefer-
1032 ence signals play complementary roles in ClaHF,
1033 and that their joint modeling is crucial for stable
1034 and efficient optimization of classification policies.

1035 **Further explanation of Figure 3.** Further anal-
1036 ysis of the interaction between α and γ reveals
1037 a clear synergistic effect. When α is small (i.e.,
1038 the RM mainly relies on relative ranking signals),
1039 varying γ yields limited gains. In contrast, when
1040 α lies in the medium range, changes in γ have
1041 a more pronounced impact on both performance
1042 and calibration, suggesting that explicit preference
1043 signals facilitate the stabilizing role of the value
1044 function. However, when $\alpha = 1$, an excessively
1045 large γ suppresses the policy’s responsiveness to
1046 preference signals and leads to performance degra-
1047 dation. This reveals a critical balance in ClaHF:
1048 preference guidance must be sufficiently strong to
1049 steer policy updates, while value constraints should
1050 stabilize training without over-regularization.

1051 Finally, we observe that the distribution of ECE
1052 does not simply trade off with Acc and F1, but
1053 reaches its minimum near the same optimal region.
1054 This indicates that ClaHF simultaneously enhances
1055 discriminative performance and confidence calibra-
1056 tion, rather than sacrificing one for the other.