

# FACTGUARD: Event-Centric and Commonsense-Guided Fake News Detection

Jing He<sup>1</sup>, Han Zhang<sup>1</sup>, Yuanhui Xiao<sup>1</sup>, Wei Guo<sup>1</sup>, Shaowen Yao<sup>1</sup>, Renyang Liu<sup>2,\*</sup>,

<sup>1</sup>School of Software and AI, Yunnan University

<sup>2</sup>Institute of Data Science, National University of Singapore

hejing@ynu.edu.cn, {hanzhang, xiaoyuanhui, guoweinyu}@stu.ynu.edu.cn, yaosw@ynu.edu.cn, ryliu@nus.edu.sg

## Abstract

Fake news detection methods based on writing style have achieved remarkable progress. However, as adversaries increasingly imitate the style of authentic news, the effectiveness of such approaches is gradually diminishing. Recent research has explored incorporating large language models (LLMs) to enhance fake news detection. Yet, despite their transformative potential, LLMs remain an untapped goldmine for fake news detection, with their real-world adoption hampered by shallow functionality exploration, ambiguous usability, and prohibitive inference costs. In this paper, we propose a novel fake news detection framework, dubbed FACTGUARD, that leverages LLMs to extract event-centric content, thereby reducing the impact of writing style on detection performance. Furthermore, our approach introduces a dynamic usability mechanism that identifies contradictions and ambiguous cases in factual reasoning, adaptively incorporating LLM advice to improve decision reliability. To ensure efficiency and practical deployment, we employ knowledge distillation to derive FACTGUARD-D, enabling the framework to operate effectively in cold-start and resource-constrained scenarios. Comprehensive experiments on two benchmark datasets demonstrate that our approach consistently outperforms existing methods in both robustness and accuracy, effectively addressing the challenges of style sensitivity and LLM usability in fake news detection.

**Code** — <https://github.com/ryliu68/FACTGUARD>

## 1 Introduction

Social media platforms have become dominant channels for information dissemination, surpassing traditional media in both scale and societal influence. However, the lack of effective content moderation on these platforms has facilitated the rapid spread of sensational fake news, further amplified by recommendation algorithms. Consequently, platforms such as Sina-Weibo and Facebook have faced significant challenges related to misinformation propagation (Refutation 2021; Avaaz 2020). Prior studies have demonstrated that the widespread circulation of fake news during major public events can trigger social panic and disrupt governance (Grinberg et al. 2019; Zhang et al. 2024; Bursztyn

et al. 2020). Given the overwhelming volume of online information, manual verification of news authenticity is infeasible, which underscores the necessity of developing automated fake news detection methods based on advanced techniques, such as deep learning, with a particular emphasis on early identification to curb the spread and societal impact of misinformation.

Early efforts in this direction have predominantly relied on insights from psychology and linguistics. In particular, a substantial body of research is rooted in psychological theories, such as the Undeutsch hypothesis (Amado, Arce, and Fariña 2015), emphasizing linguistic style differences between truthful and deceptive statements. Stylometric and emotion-based methods (Potthast et al. 2018; Rashkin et al. 2017; Ajao, Bhowmik, and Zargari 2019; Giachanou, Rosso, and Crestani 2019; Zhu et al. 2022) have thus become mainstream approaches for identifying fake news. However, as these methods primarily capture superficial features, recent studies have shown that adversaries can effectively evade detection by imitating the writing style of authentic news. This reliance on surface-level cues makes existing systems highly vulnerable to news text writing style.

To address these limitations, the research community has increasingly explored leveraging large language models (LLMs) for fake news detection. For example, some studies employ LLMs to generate adversarial samples with diverse writing styles to enhance model robustness (Wu, Guo, and Hooi 2024; Wang et al. 2024). Other approaches integrate multiple perspectives, such as combining style detection with commonsense reasoning (Hu et al. 2024), constructing multi-agent debate frameworks to aggregate different viewpoints (Liu et al. 2025) or using LLMs to simulate different news readers to generate diverse comments (Nan et al. 2024). Despite these methodological advances, several critical challenges remain unresolved. In particular, the effectiveness of style-based sample generation methods against previously unseen style attacks remains uncertain. LLM-based detection models, though promising, often suffer from low accuracy in few-shot and chain-of-thought reasoning scenarios (Hu et al. 2024), are prone to hallucinations (Xu, Jain, and Kankanhalli 2024). Besides, some methods focus only on correctly judged samples during training, while the correctness of such judgments remains unknown during inference (Hu et al. 2024), resulting in the lack of reli-

\*Corresponding author

able mechanisms for usability assessment. Moreover, multi-agent debate frameworks and role-playing-based comment generation typically incur considerable computational and time costs, which limits their practicality in cold-start<sup>1</sup> and resource-constrained environments<sup>2</sup> (Liu et al. 2025).

To bridge this gap, we propose the News Extracted Topic-Content and Commonsense Rationale Model (FACTGUARD), a comprehensive framework for robust fake news detection. FACTGUARD leverages the semantic understanding capabilities of LLMs to extract event-centric information, thereby reducing the influence of textual style. By integrating LLM-generated commonsense reasoning with advanced content extraction and dynamic reliability assessment, FACTGUARD enables a more accurate assessment of news veracity. Furthermore, the framework incorporates knowledge distillation, thereby supporting practical deployment across resource-rich<sup>3</sup>, cold-start, and resource-constrained settings and achieving a balance between accuracy and efficiency.

Specifically, in journalism communication theory, news come from the real-world events (Galtung and Ruge 1965) and style rewriting is a common deceptive strategy in fake news (Potthast et al. 2018) to accelerate news dissemination. So FACTGUARD first utilizes LLMs with carefully designed prompts to extract the core topic and principal content from news articles. This process filters out stylistic noise and preserves essential event information. To ensure the quality and relevance of the extracted content, we introduce a two-stage constraint mechanism: a text similarity metric is applied during extraction to maintain consistency with the original news, followed by an information density metric that evaluates informativeness post-extraction. The resulting event-centric content, being more objective and concise, is then semantically compared with LLM-generated commonsense rationale to enhance fake news detection. To further improve detection reliability, FACTGUARD incorporates an LLM Rationale usability module that treats the LLM as an advisor and dynamically assesses the trustworthiness of its advice via a dual-branch structure. One branch adaptively controls the influence of LLM-based judgments, while the other emphasizes potential conflicts or ambiguities identified through commonsense reasoning, enabling the model to better address complex cases. In addition, to support deployment in cost- and efficiency-sensitive scenarios, we introduce a knowledge distillation scheme (Hinton, Vinyals, and Dean 2015). This mechanism transfers knowledge from the full FACTGUARD model to a lightweight variant, FACTGUARD-D, which delivers efficient inference while preserving strong detection performance.

Extensive experiments on two widely used real-world

<sup>1</sup>Cold-start refers to the early-stage fake news detection setting where only the news content is available. In this setting, high inference efficiency is required.

<sup>2</sup>Resource-constrained refers to settings where LLMs cannot be accessed or invoked, and the model can only rely on news content for inference in such cases.

<sup>3</sup>Resource-rich refers to settings where LLMs are available and can be employed without restrictions, enabling improved fake news detection performance.

fake news detection datasets, GossipCop (Shu et al. 2020) and Weibo21 (Nan et al. 2021), demonstrate that FACTGUARD consistently outperforms state-of-the-art baselines across multiple key metrics, including accuracy and robustness. The distilled variant, FACTGUARD-D, also achieves competitive results with minimal performance degradation, confirming its practicality and robustness in resource-constrained settings. Our contributions are as follows:

- We propose FACTGUARD, a novel framework that effectively mitigates the impact of textual style on fake news detection and enables robust integration of LLM-based reasoning. To accommodate diverse practical needs, FACTGUARD is suitable for resource-rich scenarios, whereas its distilled variant, FACTGUARD-D, is tailored for cold-start and resource-constrained scenarios.
- We develop an LLM-based news extraction approach that leverages semantic understanding to obtain key topics and event content, thereby reducing style interference and enabling alignment with commonsense reasoning.
- We introduce an LLM rationale usability module, which dynamically adjusts the influence of LLM advice through a dual-branch structure based on their reliability and the presence of commonsense conflicts, ensuring the effective and adaptive use of LLM knowledge.
- We conduct extensive experiments on the GossipCop and Weibo21 datasets, which demonstrate the effectiveness and efficiency of the proposed FACTGUARD and FACTGUARD-D models.

## 2 Related Work

### 2.1 Traditional Fake News Detection

Fake news detection focuses on the early identification of misinformation, primarily based on the textual content available at publication (Qian et al. 2018). Early methods mainly relied on machine learning models with handcrafted features, including keywords, grammatical errors (Granik and Mesyura 2017), shallow linguistic patterns (Wang 2017), and statistical cues such as text length, capitalization, and punctuation (Castillo, Mendoza, and Poblete 2011). With advances in deep learning, LSTM-based approaches were introduced to capture linguistic differences in news with satirical or rumor styles (Rashkin et al. 2017), and some studies explored sentiment-related features (Ajao, Bhowmik, and Zargari 2019; Giachanou, Rosso, and Crestani 2019). However, these methods fundamentally rely on surface-level features, making them vulnerable to variations in writing style.

To address these limitations, style modeling with pre-trained language models such as BERT and RoBERTa (Przybyla 2020) has become common in fake news detection. Nevertheless, approaches relying solely on textual features remain inadequate for combating increasingly sophisticated misinformation. Recent studies have therefore incorporated background knowledge beyond textual content, including social context (Shu et al. 2019; Cui et al. 2022), social emotion (Zhang et al. 2021), news environment (Sheng et al. 2022), and external knowledge (Hu et al. 2022). While small

language models (SLMs) offer certain improvements, their limited knowledge and capacity continue to constrain further progress in fake news detection.

## 2.2 LLM-based Fake News Detection

Recent studies have leveraged the strong language generation and comprehension capabilities of LLMs for content generation and enhancement in fake news detection. For content generation, LLM-Fake analyzes the psychological motivations behind LLM-generated fake news and constructs the MegaFake dataset (Wang et al. 2024); SheepDog employs LLM-generated multi-style samples as adversarial data to improve detector robustness (Wu, Guo, and Hooi 2024). For content enhancement, ARG leverages prompt engineering to guide LLM in multi-perspective analysis, with SLM integrating the final judgment (Hu et al. 2024). LEKD incorporates offline LLM knowledge via semantic graph alignment and knowledge distillation (Chen et al. 2025). TED introduces a structured multi-agent debate mechanism, enabling LLM to reason from diverse perspectives (Liu et al. 2025). GenFEND simulate readers with different identities to generate diverse comments, thereby providing additional information for early-stage detection (Nan et al. 2024).

Despite these advances, the fundamental problem of style sensitivity remains unresolved. While LLM-based content generation increases sample diversity, it often fails to capture event semantics and does not fully eliminate style-related interference (Wang et al. 2024; Wu, Guo, and Hooi 2024). In addition, augmented data generated by LLMs is often not effectively integrated with the detection backbone, resulting in limited overall improvement (Hu et al. 2024). Excessive reliance on LLMs and their high inference costs further constrain practical deployment (Hu et al. 2024; Liu et al. 2025; Nan et al. 2024), particularly in cold-start or resource-constrained scenarios.

In summary, vulnerability to writing style, insufficient event-centric modeling, and the challenge of efficiently integrating LLM capabilities remain open problems. Motivated by these gaps, this paper proposes a new framework, FACTGUARD, that systematically addresses style interference, leverages LLMs’ commonsense reasoning capacity, and supports efficient deployment for fake news detection.

## 3 Framework

### 3.1 Preliminary

**Problem Formulation.** In resource-rich scenarios, we consider a dataset  $\mathcal{D}_{\text{news}} = \{(n_i, c_i, r_i)\}_{i=1}^N$ , where  $n_i$  denotes the news text, while  $c_i$  and  $r_i$  represent the event-based topic-content and the commonsense rationale extracted by an LLM. The goal for each news item  $D_i = (n_i, c_i, r_i)$  is to predict  $\hat{y}_i \in \{0, 1\}$ , where 1 indicates fake news and 0 indicates true news, by integrating the original text with LLM-generated features. The detection model (i.e., FACTGUARD) encodes  $n_i$  via a dual-attention module  $f_\theta(n_i)$ , fuses  $c_i$  and  $r_i$  with cross-attention and usability assessment  $g_\phi(\text{CA}(c_i, r_i))$ , and concatenates both for final prediction:

$$\hat{y}_i = \text{MLP}([f_\theta(n_i); g_\phi(\text{CA}(c_i, r_i))]) \quad (1)$$

where  $[\ ]$  denotes vector concatenation. To further address cold-start and resource-constrained scenarios, we distill an efficiency and resource-friendly student model (i.e., FACTGUARD-D) from the teacher model FACTGUARD, where the FACTGUARD-D only takes  $n_i$  as input and predicts its label  $\hat{y}_i$ :

$$\hat{y}_i = f_\psi(n_i) \quad (2)$$

**Notations.** All notations used throughout this paper are summarized in Table 3 in Appendix B.

### 3.2 FACTGUARD Overview

Figure 1(a–c) illustrates the main modules as well as the training and inference procedures of FACTGUARD. The core goal of FACTGUARD is to achieve style debiasing and fully exploit the capabilities of LLMs for robust fake news detection. For each news item  $n$ , the model first employs the LLM to extract topic-content information  $c$  as well as commonsense rationale  $r$ . These elements, together with the original news text, are encoded using an SLM (see Figure 1(a)). The Topic-Content&Rationale Interactor enables deep feature interaction between the extracted topic-content and the commonsense rationale, while the Rationale Usability Evaluator adaptively assigns weights to the LLM-provided advice. The resultant interacted features  $f_{llm}$  are then aggregated with the news features  $f_N$ . The fused representations are utilized for veracity prediction for  $n$  (see Figure 1(b)). During training, three loss functions— $L_{cls}$ ,  $L_{usability}$ , and  $L_{text}$ —are employed to optimize model parameters. Once FACTGUARD is well trained, it can be used to predict the veracity of the unseen news sample  $n$  by leveraging the inputs  $c$ ,  $r$ , and  $n$  (see Figure 1(c)).

### 3.3 Feature Extraction

Pretrained SLMs such as BERT or RoBERTa are trained on large-scale datasets in an unsupervised fashion, enabling them to generate high-dimensional contextual representations well-suited for various downstream tasks. To effectively extract information features, we employ these models as text encoders within our framework. Specifically, for a given news item  $n$ , the extracted topic-content  $c$ , and the commonsense rationale  $r$ , we denoted them as  $N$  (news),  $C$  (topic-content), and  $R$  (commonsense rationale), respectively.

### 3.4 Feature Concatenation

This module aims to obtain high-quality representations for both the LLM-generated augmented information and the original news content, facilitating their effective integration and collaboration as the foundation for fake news detection. Feature integration is performed via concatenation of the respective representations, enabling subsequent modules to leverage comprehensive contextual information.

**Topic-Content&Rationale Interactor.** To enable comprehensive feature exchange between the LLM-extracted topic-content and commonsense rationale, we introduce a dual cross-attention module based on multi-head attention.

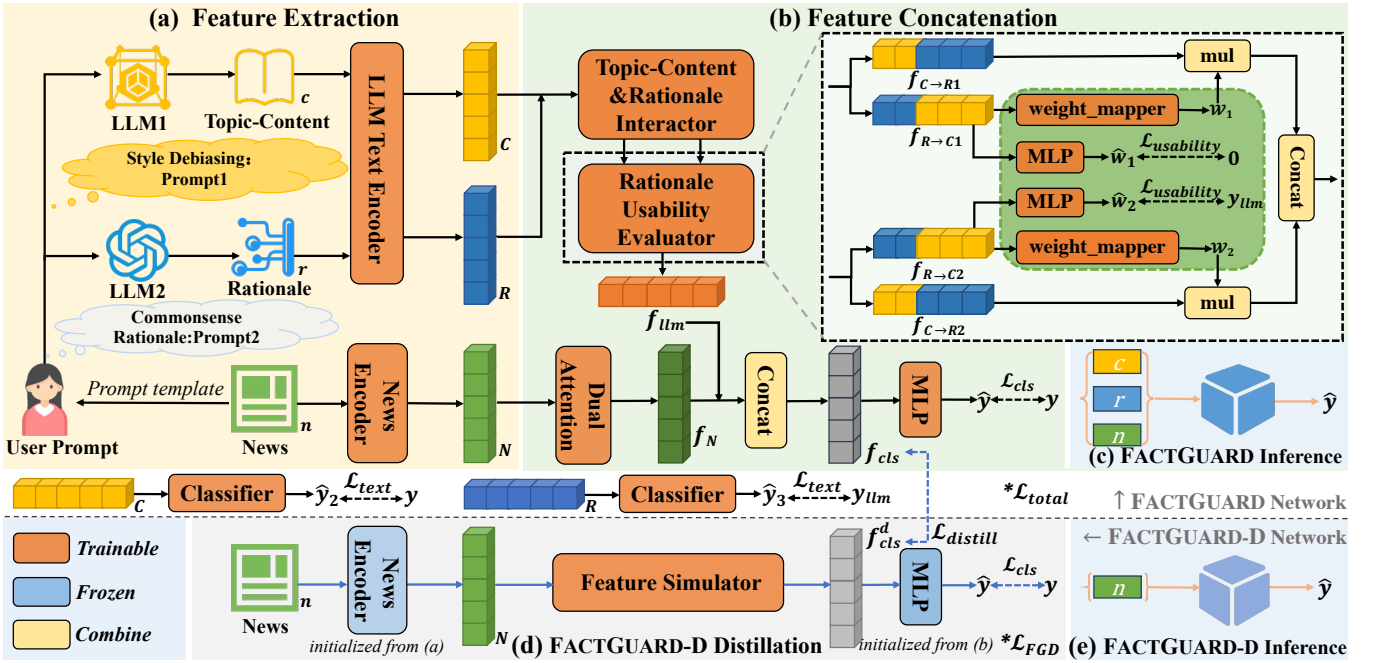


Figure 1: Overview of FACTGUARD and FACTGUARD-D. FACTGUARD main consists of two modules: (1) **Feature Extraction**, which identifies topic content and enables commonsense reasoning for each news article using an LLM. The resulting features and the original text are encoded for downstream processing. (2) **Feature Concatenation**, which adaptively integrates LLM-derived features with news content via a cross-attention mechanism and the Rationale Usability Evaluator, followed by MLP-based classification. After training, knowledge distillation yields a lightweight FACTGUARD-D without LLMs’ advice.

The computation is formulated as follows:

$$\text{head}_i = \text{softmax} \left( \frac{Q_i K_i^\top}{\sqrt{d_k}} \right) V_i, \quad (3)$$

$$\text{CA}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O, \quad (4)$$

where  $Q_i = QW_i^Q$ ,  $K_i = KW_i^K$ , and  $V_i = VW_i^V$ . Here,  $d_k$  is the dimension of each attention head,  $h$  is the number of heads, and  $W^O$  is the output projection matrix. Given topic-content  $c$  and commonsense rationale  $r$ , after embedding, the interactions are computed as:

$$f_{C \rightarrow R} = \text{AvgPool}(\text{CA}(C, R, R)), \quad (5)$$

$$f_{R \rightarrow C} = \text{AvgPool}(\text{CA}(R, C, C)), \quad (6)$$

where  $\text{AvgPool}(\cdot)$  denotes average pooling applied over the token representations output by the cross-attention layer.  $f_{C \rightarrow R}$  represents the LLM advice feature vector, and  $f_{R \rightarrow C}$  serves as a weighting factor in the Rationale Usability Evaluator.

**Rationale Usability Evaluator.** Directly enforcing consistency between LLM judgments and ground-truth labels may result in the loss of valuable complementary features provided by the LLMs. To address this, we propose a rationale usability evaluation module to dynamically adjust the fusion weights of LLM features, thereby maximizing the utility of LLM-generated knowledge. This module adopts a dual-branch MLP structure: one branch reduces its contribution as the LLMs’ direct detection capability is limited,

while the other increases its contribution as commonsense reasoning can provide more effective information when it identifies contradictions or uncertainty. A three-layer MLP (weight\_mapper) maps the feature vectors  $f_{R \rightarrow C_i}$  to fusion weights  $w_i$  as follows:

$$w_i = \text{sigmoid}(\text{weight\_mapper}(f_{R \rightarrow C_i})), \quad i = 1, 2. \quad (7)$$

The final LLM feature representation  $f_{llm}$  is then computed as:

$$f_{llm} = [w_1 \cdot f_{C \rightarrow R1}; w_2 \cdot f_{C \rightarrow R2}], \quad (8)$$

where  $w_1$  and  $w_2$  are the fusion weights for the two branches, and  $f_{C \rightarrow R1}$ ,  $f_{C \rightarrow R2}$  denote their respective interaction features.

**Dual Attention Fusion.** To further improve the expressiveness and robustness of features derived from BERT or RoBERTa, we introduce a linear attention mechanism that adaptively assigns higher weights to salient tokens, thereby suppressing irrelevant or noisy information:

$$\text{Attn}(X) = \sum_{t=1}^T \text{softmax}(Wx_t + b) \cdot x_t, \quad (9)$$

where  $x_t$  denotes the input feature at position  $t$ ,  $W$  and  $b$  are learnable parameters, and  $T$  is the token sequence length. To enhance model robustness, a dual-branch architecture is adopted, where the same linear attention module is applied in parallel to both branches. The outputs of the two branches

are then averaged to obtain the final news feature representation:

$$f_N = \frac{\text{Attn}(N) + \text{Attn}(N)}{2}. \quad (10)$$

**Feature Concatenation.** Based on the outputs obtained in the previous step, the news feature vector  $f_N$  and the LLM-enhanced feature vector  $f_{llm}$  are summed to facilitate the final prediction. For each news item  $n$  with label  $y \in \{0, 1\}$ , these vectors are combined to produce the final feature representation  $f_{cls}$ , computed as:

$$f_{cls} = [f_N; f_{llm}]. \quad (11)$$

$f_{cls}$  is subsequently input into an MLP classifier to predict the veracity label:

$$\hat{y} = \text{MLP}(f_{cls}). \quad (12)$$

### 3.5 Training

**Data Process.** To enhance semantic understanding and reduce the influence of writing style, LLMs is leveraged to extract the topic-content of each news article. Additionally, commonsense rationale analysis is performed by the LLMs to identify and judge content that may contradict commonsense. Due to page limitations, detailed information on the data processing procedure is provided in Appendix D.

**Objective.** The FACTGUARD method is designed with three principal objectives: ① to achieve accurate prediction of news veracity; ② to effectively integrate model recommendations and fully leverage the capabilities of LLMs; and ③ to enhance the representation of information augmentation provided by LLMs. Accordingly, the overall objective loss function is defined as a weighted sum of the prediction loss, the LLM rationale usability loss, and the information augmentation representation loss.

To improve the final detection performance, the Binary Cross-Entropy (BCE) classification loss is computed to guide the model in accurately identifying fake news:

$$\mathcal{L}_{cls} = \text{BCE}(\hat{y}, y). \quad (13)$$

To supervise the learning of LLM features, the supervision signals for the weights are set as 0 for one branch and  $y_{llm}$  for the other:

$$\hat{w}_i = \text{sigmoid}(\text{MLP}(f_{C \rightarrow R_i})), \quad i = 1, 2, \quad (14)$$

$$\mathcal{L}_{usability} = \text{BCE}(\hat{w}_1, 0) + \text{BCE}(\hat{w}_2, y_{llm}). \quad (15)$$

To enhance the utility of LLM-generated augmentations, we employ an auxiliary task that aligns extracted semantics and commonsense reasoning with ground-truth labels and LLM veracity judgments. Augmented representations are fed into a classifier composed of a linear attention and an MLP head (without sigmoid), optimized by cross-entropy (CE) loss:

$$\hat{y}_2 = \text{Classifier}(C), \quad \hat{y}_3 = \text{Classifier}(R), \quad (16)$$

$$\mathcal{L}_{text} = \text{CE}(\hat{y}_2, y) + \text{CE}(\hat{y}_3, y_{llm}). \quad (17)$$

The total loss function is a weighted sum of the aforementioned terms:

$$\mathcal{L}_{total} = \mathcal{L}_{cls} + \alpha \frac{\mathcal{L}_{usability}}{2} + \beta \frac{\mathcal{L}_{text}}{2}, \quad (18)$$

where  $\alpha$  and  $\beta$  are hyperparameter weights, and the division by 2 is used to average the two sub-losses, ensuring balanced contributions from each component in the overall loss.

### 3.6 Inference

In resource-rich scenarios, LLMs are leveraged via prompt engineering to extract the topic-content, and commonsense rationales of news articles. The extracted outputs, together with the original news text, are subsequently fed into the well-trained frozen FACTGUARD model for veracity prediction.

Due to page limitations, further training and inference details for FACTGUARD are provided in Algorithm 1 and Algorithm 2 in Appendix C.

### 3.7 FACTGUARD-D

**Distillation.** Directly invoking LLMs for each prediction in FACTGUARD is impractical in resource-constrained or latency-sensitive cold-start scenarios due to the substantial overhead of real-time LLM prompting for text extraction and commonsense reasoning. To address this, we develop a llm-free student model via knowledge distillation from the trained FACTGUARD, following a teacher-student paradigm (Hu et al. 2024). The core idea is to transfer and internalize the teacher model’s reasoning knowledge into a parameterized lightweight student network. Specifically, as illustrated in Figure 1(d), the student model’s news encoder and classifier are initialized from the trained FACTGUARD. To acquire the teacher’s reasoning capabilities, a feature simulator is implemented as a four-layer Transformer encoder and a linear attention module to internalize the teacher’s knowledge. In addition to the standard prediction loss  $\mathcal{L}_{cls}$  as in FACTGUARD, the student model is further supervised by a feature distillation loss  $\mathcal{L}_{distill}$ , which encourages the student’s feature representation  $f_{cls}^d$  to approximate that of the teacher  $f_{cls}$  by minimizing the mean squared error (MSE) between them:

$$\mathcal{L}_{distill} = \text{MSE}(f_{cls}^d, f_{cls}). \quad (19)$$

**Inference.** In cold-start and resource-constrained scenarios, the FACTGUARD-D model operates exclusively on the original news text  $n$ , achieving fast predictions with only a slight reduction in accuracy.

Due to page limitations, further distillation training and inference details for FACTGUARD-D are provided in Algorithm 3 and Algorithm 4 in Appendix C.

## 4 Experiments

This section presents comprehensive experimental studies of the proposed FACTGUARD and FACTGUARD-D models. We first introduce the experimental setup (for detailed settings, please refer to Appendix E). Subsequently, we compare FACTGUARD with a wide range of baselines, conduct ablation studies to assess the contribution of each model component, analyze parameter sensitivity, and discuss challenges associated with LLM-based text extraction.

Group	Model	Weibo21				GossipCop			
		macF1	Acc.	F1 <sub>real</sub>	F1 <sub>fake</sub>	macF1	Acc.	F1 <sub>real</sub>	F1 <sub>fake</sub>
G1	GPT-3.5-turbo*	0.725	0.734	0.774	0.676	0.702	0.813	0.884	0.519
	GPT-4o-mini#	0.725	0.746	0.780	0.670	0.691	0.845	0.909	0.472
	ChatEval-o#	0.694	0.717	0.778	0.611	0.733	0.860	0.919	0.546
	ChatEval-s#	0.694	0.719	0.780	0.608	0.738	0.869	0.923	0.553
G2	BERT*	0.753	0.754	0.769	0.737	0.765	0.862	0.916	0.615
	RoBERTa	0.753	0.755	0.775	0.731	0.765	0.862	0.916	0.613
	EANN*	0.754	0.756	0.773	0.736	0.763	0.864	0.918	0.608
	Publisher-Emo*	0.761	0.763	0.784	0.738	0.766	0.868	0.920	0.611
	ENDEF*	0.765	0.766	0.779	0.751	0.768	0.865	0.918	0.618
G3	Bert + Rationale*	0.767	0.769	0.787	0.748	0.777	0.870	0.921	0.633
	SuperICL*	0.757	0.759	0.779	0.734	0.736	0.864	0.920	0.551
	Bert + GenFEND	0.755	0.760	0.791	0.719	0.764	0.875	0.926	0.603
	Roberta + GenFEND	0.771	0.774	0.796	0.747	0.770	0.866	0.919	0.621
	ARG*	0.784	0.786	0.804	0.764	0.790	0.878	0.926	0.653
	TED#	<u>0.795</u>	<u>0.798</u>	<u>0.815</u>	<u>0.774</u>	<u>0.803</u>	<b>0.892</b>	0.932	<u>0.674</u>
	<b>Ours</b>	<b>0.801</b>	<b>0.804</b>	<b>0.824</b>	<b>0.777</b>	<b>0.805</b>	<b>0.892</b>	<b>0.935</b>	<b>0.675</b>
G4	ARG-D*	0.771	0.772	0.785	0.756	0.778	0.870	0.921	0.634
	<b>Ours</b>	0.788	0.790	0.807	0.769	0.790	0.888	<u>0.933</u>	0.647

Table 1: Performance comparison on Weibo21 and GossipCop datasets across four metrics, i.e., macF1, Accuracy, F1<sub>real</sub>, and F1<sub>fake</sub>. The highest result in each category is **bolded** and the second highest result is underlined. In the results table, \* means the result is from (Hu et al. 2024) and # means the result is from (Liu et al. 2025).

## 4.1 Setup

**Datasets.** We employ the Weibo21 (Chinese) (Nan et al. 2021) and GossipCop (English) (Shu et al. 2020) for evaluation. Both datasets are preprocessed by deduplication and temporal splitting, following established practices (Zhu et al. 2022; Mu, Bontcheva, and Aletras 2023; Hu et al. 2024), to mitigate the risk of data leakage and prevent overestimation of SLM performance. In addition, we also utilize the commonsense rationales from (Hu et al. 2024).

**Baselines.** Recent fake news detection methods predominantly rely on LLMs and SLMs, and can be categorized into four groups. Among them, we involved 14 representative baselines in this work. The first group (G1) comprises LLM-only methods, including GPT-3.5-turbo (OpenAI 2023), GPT-4o-mini (OpenAI 2024), ChatEval-o (one-by-one strategy) (Chan et al. 2024), and ChatEval-s (Simultaneous-Talk strategy) (Chan et al. 2024). The second group (G2) consists of SLM-only methods, such as BERT (Devlin et al. 2019), RoBERTa (Liu et al. 2019), EANN (Wang et al. 2018), Publisher-Emo (Zhang et al. 2021), and ENDEF (Zhu et al. 2022). The third group (G3) includes LLM-SLM methods, such as BERT + Rationale (Hu et al. 2024), SuperICL (Zhong et al. 2023), ARG (Hu et al. 2024), BERT + GenFEND (Nan et al. 2024), RoBERTa + GenFEND and TED (Liu et al. 2025). The fourth group (G4) comprises methods employing model distillation, such as ARG-D (Hu et al. 2024).

**Metrics.** We evaluate performance using four metrics: Accuracy (Acc.), F1<sub>real</sub>, F1<sub>fake</sub>, and Macro-F1 (macF1).

**Implementation Details.** We utilize bert-base-chinese<sup>4</sup> (Devlin et al. 2019) as the text encoder for Chinese FACTGUARD model and roberta-base<sup>5</sup> (Liu et al. 2019) for English FACTGUARD model. For the Weibo21 and GossipCop dataset, topic-content extraction is performed using locally deployed DeepSeek-R1-Distill-Llama-8B<sup>6</sup> (DeepSeek-AI 2025) and SOLAR-10.7B-Instruct-v1.0-uncensored<sup>7</sup> (Kim et al. 2024) (Kim et al. 2024), respectively. Commonsense reasoning modules for both datasets are adopted from prior work (Hu et al. 2024). We employ the AdamW optimizer with an initial learning rate of  $2e-4$  and a weight decay of  $5e-5$ . Early stopping with a patience of 5 epochs is applied to prevent overfitting. All experiments are conducted on a single NVIDIA A100 (40GB) GPU with a fixed random seed of 3759, PyTorch version 1.13.0.

## 4.2 Comparative Results

To evaluate the effectiveness of the proposed FACTGUARD model, we conduct systematic experiments on the Weibo21 and GossipCop datasets. As shown in Table 1, FACTGUARD consistently outperforms all baseline methods across multiple evaluation metrics, achieving the best results across both datasets. These results underscore the superior cross-lingual generalization and stability of the proposed model.

<sup>4</sup><https://huggingface.co/google-bert/bert-base-chinese>

<sup>5</sup><https://huggingface.co/FacebookAI/roberta-base>

<sup>6</sup><https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B>

<sup>7</sup><https://huggingface.co/upstage/SOLAR-10.7B-Instruct-v1.0>

Group	Model	Weibo21				GossipCop			
		macF1	Acc.	F1 <sub>real</sub>	F1 <sub>fake</sub>	macF1	Acc.	F1 <sub>real</sub>	F1 <sub>fake</sub>
G1	<b>FACTGUARD</b>	<b>0.801</b>	<b>0.804</b>	<b>0.824</b>	<b>0.777</b>	<b>0.805</b>	<b>0.892</b>	<b>0.935</b>	<b>0.675</b>
G2	News	0.768	0.769	0.784	0.751	0.765	0.862	0.916	0.613
	Topic-Content	0.690	0.691	0.708	0.672	0.769	0.861	0.915	0.624
	Commonsense	0.678	0.685	0.728	0.627	0.698	0.832	0.899	0.498
G3	w/o News	0.718	0.722	0.753	0.683	0.773	0.873	0.924	0.623
	w/o Topic-Content	0.772	0.774	0.794	0.750	0.779	0.875	0.925	0.632
	w/o Commonsense	0.770	0.773	0.797	0.743	0.779	0.871	0.922	0.633
G4	w/o llm-usability	0.778	0.780	0.794	0.763	0.780	0.878	0.927	0.633
	use ARG-usefulness	0.782	0.782	0.793	0.770	0.775	0.872	0.923	0.628

Table 2: An ablation study of FACTGUARD on the Weibo21 and GossipCop datasets evaluates the contribution of each model component. Specifically, G2 assesses the predictive power of individual input features (original news, LLM-extracted topic-content, LLM commonsense rationale judgment); G3 quantifies the impact of ablating each core input module; and G4 investigates the effects of removing LLM rationale usability module on overall performance and cross-lingual generalization.

**Weibo21.** On the Weibo21 dataset, FACTGUARD outperforms the strongest baseline TED by 0.8% in accuracy and 0.9% in F1<sub>real</sub>, along with notable macF1 gains. These improvements arise from its LLM-driven topic-content extraction, dual-branch reasoning, and efficient rationale usability module, which jointly combine SLM efficiency with LLM reasoning capacity. Compared with TED’s complex multi-agent debate framework, FACTGUARD achieves higher accuracy with only two simple prompts, significantly reducing computational resource and inference costs. The distilled variant, FACTGUARD-D, further surpasses ARG and ARG-D, demonstrating effective compression and strong practical performance.

**GossipCop.** On the GossipCop dataset, although the improvements are smaller, FACTGUARD still attains the best macF1 (0.805) and F1<sub>real</sub> (0.935). Under limited computational resources, FACTGUARD-D also achieves lower misclassification rates, confirming its broad applicability. The performance differences across datasets mainly result from class imbalance (more fake news in Weibo21, more real news in GossipCop), yet FACTGUARD remains stable and generalizable across languages.

In summary, both FACTGUARD and FACTGUARD-D demonstrate strong performance, cross-lingual robustness, and model compression capability, making them efficient and reliable solutions for multilingual fake news detection.

### 4.3 Ablation Study

We conduct ablation studies on the GossipCop and Weibo21 datasets to evaluate the effectiveness of each module in the proposed FACTGUARD framework, focusing on the LLM-based topic-content extraction, commonsense rationale judgment, and usability evaluation modules. The ablation experiments are organized into three groups:

As shown in Table 2, the main findings are as follows: (1) Removing the original news representation yields the largest reduction in overall performance, confirming its indispensable role as the foundation for fake news detection. (2) Removing the LLM-based topic-content extraction module re-

sults in a notable drop in macF1 and F1<sub>real</sub>, highlighting its importance in capturing essential event information and reducing stylistic noise. (3) Excluding the LLM commonsense rationale module further degrades performance, demonstrating its value in improving factual consistency and reasoning. (4) Omitting the usability evaluation module or use ARG’s LLM usefulness module also leads to decreased performance, underscoring its role in aligning LLM inference with robust news representations. (5) The LLM-extracted topic-content module and the LLM commonsense rationale module need to be used in conjunction to achieve the maximum performance improvement. Overall, these results validate that each component is essential for FACTGUARD’s strong performance in multilingual fake news detection.

In addition to the main ablation experiments, we performed both grid search and random search to determine optimal loss weight parameters for the multi-objective loss functions in the FACTGUARD model. For the Weibo21 configuration, the best results were achieved with  $\alpha = 0.40$  and  $\beta = 0.16$ , while for the GossipCop configuration, the optimal values were  $\alpha = 0.50$  and  $\beta = 0.58$ . To facilitate effective distillation, we set the distillation coefficient  $\lambda$  in the loss function to 8 for both the Chinese and English models in FACTGUARD-D. Additional experimental details are provided in Appendix F.

## 5 Conclusion

We propose FACTGUARD, a model that leverages LLMs for semantic understanding and commonsense reasoning to improve fake news detection performance. By extracting topic and core content and employing a usability evaluation module in commonsense rationale, FACTGUARD effectively reduces style bias and integrates LLM-generated judgments. For cold-start and resource-limited scenarios, the distilled variant FACTGUARD-D is optimized for efficiency and resources. Experiment results on Weibo21 and GossipCop datasets show that FACTGUARD outperforms baselines with each module proven effective by ablation studies, while FACTGUARD-D achieves a strong balance between accuracy

and speed.

**Future Work.** Future directions include: (1) developing customized methods for Chinese and English fake news; (2) optimizing the model for edge deployment; (3) enhancing interpretability of usability evaluation module to improve transparency and credibility; (4) exploring the role of text style at different stages of news dissemination detection; and (5) considering the benchmark data contamination of the employed LLMs, and extending cross-domain adaptation across emerging platforms and multimodal signals.

## 6 Acknowledgments

This work was supported in part by Scientific Research and Innovation Project of Postgraduate Students in the Academic Degree of Yunnan University under KC-252512080, in part by the National Natural Science Foundation of China under Grants 62162067 and 82360280, in part by the Yunnan Province Special Project under Grant 202403AP140021 and in part by the Yunnan Fundamental Research Project under Grant 202401AT070474.

## References

- Ajao, O.; Bhowmik, D.; and Zargari, S. 2019. Sentiment aware fake news detection on online social networks. In *ICASSP*, 2507–2511.
- Amado, B. G.; Arce, R.; and Fariña, F. 2015. Undeutsch hypothesis and Criteria Based Content Analysis: A meta-analytic review. *The European Journal of Psychology Applied to Legal Context*, 7: 3–12.
- Avaaz. 2020. Facebook’s algorithm: A major threat to public health. [https://secure.avaaz.org/campaign/en/facebook-threat\\_health/](https://secure.avaaz.org/campaign/en/facebook-threat_health/). Accessed: 2025-07-22.
- Bursztyn, L.; Rao, A.; Roth, C. P.; and Yanagizawa-Drott, D. H. 2020. Misinformation During a Pandemic. Technical Report 27417, National Bureau of Economic Research.
- Cao, Z.; Nguyen, J.; and Zafarani, R. 2025. Is Less Really More? Fake News Detection with Limited Information. In *KDD*.
- Castillo, C.; Mendoza, M.; and Poblete, B. 2011. Information credibility on twitter. In *WWW*, 675–684.
- Chan, C.-M.; Chen, W.; Su, Y.; Yu, J.; Xue, W.; Zhang, S.; Fu, J.; and Liu, Z. 2024. ChatEval: Towards Better LLM-based Evaluators through Multi-Agent Debate. In *ICLR*.
- Chen, X.; Huang, X.; Gao, Q.; Huang, L.; and Liu, G. 2025. Enhancing text-centric fake news detection via external knowledge distillation from LLMs. *Neural Networks*, 187: 107377.
- Cui, J.; Kim, K.; Na, S. H.; and Shin, S. 2022. Meta-path-based fake news detection leveraging multi-level social context information. In *CIKM*, 325–334.
- Cui, Y.; Che, W.; Liu, T.; Qin, B.; Wang, S.; and Hu, G. 2020. Revisiting Pre-Trained Models for Chinese Natural Language Processing. In *EMNLP*, 657–668.
- DeepSeek-AI. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. arXiv:2501.12948.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 4171–4186.
- Galtung, J.; and Ruge, M. H. 1965. The structure of foreign news: The presentation of the Congo, Cuba and Cyprus crises in four Norwegian newspapers. *Journal of peace research*, 2: 64–90.
- Giachanou, A.; Rosso, P.; and Crestani, F. 2019. Leveraging emotional signals for credibility detection. In *SIGIR*, 877–880.
- Granik, M.; and Mesyura, V. 2017. Fake news detection using naive Bayes classifier. In *IEEE UkrCon*, 900–903.
- Grinberg, N.; Joseph, K.; Friedland, L.; Swire-Thompson, B.; and Lazer, D. 2019. Fake news on Twitter during the 2016 US presidential election. *Science*, 363: 374–378.
- Hinton, G.; Vinyals, O.; and Dean, J. 2015. Distilling the knowledge in a neural network. arXiv:1503.02531.
- Hu, B.; Sheng, Q.; Cao, J.; Shi, Y.; Li, Y.; Wang, D.; and Qi, P. 2024. Bad actor, good advisor: Exploring the role of large language models in fake news detection. In *AAAI*, 22105–22113.
- Hu, X.; Guo, Z.; Wu, G.; Liu, A.; Wen, L.; and Yu, P. S. 2022. CHEF: A pilot Chinese dataset for evidence-based fact-checking. In *NAACL*, 3362–3376.
- Kim, S.; Kim, D.; Park, C.; Lee, W.; Song, W.; Kim, Y.; Kim, H.; Kim, Y.; Lee, H.; Kim, J.; et al. 2024. SOLAR 10.7 B: Scaling Large Language Models with Simple yet Effective Depth Up-Scaling. In *NAACL*, 23–35.
- Liu, Y.; Liu, Y.; Zhang, X.; Chen, X.; and Yan, R. 2025. The truth becomes clearer through debate! Multi-agent systems with large language models unmask fake news. In *SIGIR*, 504–514.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv:1907.11692.
- Mu, Y.; Bontcheva, K.; and Aletras, N. 2023. It’s about Time: Rethinking Evaluation on Rumor Detection Benchmarks using Chronological Splits. In *EACL*, 736–743.
- Nan, Q.; Cao, J.; Zhu, Y.; Wang, Y.; and Li, J. 2021. MD-FEND: Multi-domain Fake News Detection. In *CIKM*, 3343–3347.
- Nan, Q.; Sheng, Q.; Cao, J.; Hu, B.; Wang, D.; and Li, J. 2024. Let silence speak: Enhancing fake news detection with generated comments from large language models. In *CIKM*, 1732–1742.
- OpenAI. 2023. GPT-3.5 Turbo Model. <https://platform.openai.com/docs/models/gpt-3.5-turbo>. Accessed: 2025-07-22.
- OpenAI. 2024. GPT-4o mini: Advancing Cost-Efficient Intelligence. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>. Accessed: 2025-07-22.
- Potthast, M.; Kiesel, J.; Reinartz, K.; Bevendorff, J.; and Stein, B. 2018. A Stylometric Inquiry into Hyperpartisan



and Fake News. In *ACL*, 231–240. Melbourne, Australia: Association for Computational Linguistics.

Przybyla, P. 2020. Capturing the style of fake news. In *AAAI*, 490–497. New York, USA: AAAI Press.

Qian, F.; Gong, C.; Sharma, K.; and Liu, Y. 2018. Neural User Response Generator: Fake News Detection with Collective User Intelligence. In *IJCAI*, 3834–3840.

Rashkin, H.; Choi, E.; Jang, J. Y.; Volkova, S.; and Choi, Y. 2017. Truth of Varying Shades: Analyzing Language in Fake News and Political Fact-Checking. In *EMNLP*, 2931–2937. Copenhagen, Denmark: Association for Computational Linguistics.

Refutation, W. R. 2021. Weibo Rumor Refutation Report 2021. <https://weibo.com/1866405545/LcFuud7ml#repost>. Accessed: 2025-07-22.

Sheng, Q.; Cao, J.; Zhang, X.; Li, R.; Wang, D.; and Zhu, Y. 2022. Zoom Out and Observe: News Environment Perception for Fake News Detection. In *ACL*, 4543–4556.

Shu, K.; Cui, L.; Wang, S.; Lee, D.; and Liu, H. 2019. defend: Explainable fake news detection. In *KDD*, 395–405.

Shu, K.; Mahudeswaran, D.; Wang, S.; Lee, D.; and Liu, H. 2020. Fakenewsnet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media. *Big data*, 8: 171–188.

Team, Q. 2025. Qwen3 Technical Report. arXiv:2505.09388.

Wang, L. Z.; Ma, Y.; Gao, R.; Guo, B.; Zhu, H.; Fan, W.; Lu, Z.; and Ng, K. C. 2024. Megafake: a theory-driven dataset of fake news generated by large language models. arXiv:2408.11871.

Wang, W. Y. 2017. “Liar, Liar Pants on Fire”: A New Benchmark Dataset for Fake News Detection. In *ACL*, 422–426.

Wang, Y.; Ma, F.; Jin, Z.; Yuan, Y.; Xun, G.; Jha, K.; Su, L.; and Gao, J. 2018. Eann: Event adversarial neural networks for multi-modal fake news detection. In *KDD*, 849–857.

Wu, J.; Guo, J.; and Hooi, B. 2024. Fake News in Sheep’s Clothing: Robust Fake News Detection Against LLM-Empowered Style Attacks. In *KDD*, 3367–3378.

Xu, Z.; Jain, S.; and Kankanhalli, M. 2024. Hallucination is inevitable: An innate limitation of large language models. arXiv:2401.11817.

Zhang, X.; Cao, J.; Li, X.; Sheng, Q.; Zhong, L.; and Shu, K. 2021. Mining dual emotion for fake news detection. In *WWW*, 3465–3476.

Zhang, Z.; Liu, Q.; Hu, Z.; Zhan, Y.; Huang, Z.; Gao, W.; and Mao, Q. 2024. Enhancing Fairness in Meta-learned User Modeling via Adaptive Sampling. In *WWW*, 3241–3252. Singapore: ACM.

Zhong, Q.; Ding, L.; Liu, J.; Du, B.; and Tao, D. 2023. Can chatgpt understand too? a comparative study on chatgpt and fine-tuned bert. arXiv:2302.10198.

Zhu, Y.; Sheng, Q.; Cao, J.; Li, S.; Wang, D.; and Zhuang, F. 2022. Generalizing to the future: Mitigating entity bias in fake news detection. In *SIGIR*, 2120–2125.

## A Appendix Overview

This appendix provides supplementary material that could not be included in the main paper due to space constraints. Specifically, it includes:

- **Section. B:** Introduction of all notations used throughout the paper.
- **Section. C:** Complete algorithmic procedures for training and inference in the proposed FACTGUARD and FACTGUARD-D frameworks.
- **Section. D:** LLM-based news enhancement, including topic and content extraction, extraction results, and commonsense rationale-based news judgment.
- **Section. E:** Details of datasets, baselines and evaluation metrics used in comparative experiments.
- **Section. F:** Additional ablation studies on loss hyperparameters sensitivity of FACTGUARD and FACTGUARD-D, text encoders’ choices, confidence distribution analysis and case analysis.

## B Notation

Notations used in this paper are summarized in Table 3.

Notation	Description
$c$	LLM-extracted news topic-content
$r$	LLM’s commonsense rationale on news
$n$	Original news content
$C$	Topic-Content’s token
$R$	Commonsense rationale’s token
$N$	News’ token
$f_{C \rightarrow R}$	LLM advice features
$f_{R \rightarrow C}$	Used as the weight of $f_{C \rightarrow R}$ after score map
$w_i$	Adjusted weight of LLM advice
$\hat{w}_i$	Predicted weight of LLM advice
$f_{llm}$	The final fused feature of $c$ and $r$
$f_N$	The final feature of $n$
$f_{cls}$	The final feature for prediction
$y$	Label of news
$y_{llm}$	Label of LLM’s judgment
$\hat{y}_2$	Prediction of $c$
$\hat{y}_3$	Prediction of $r$
$L_{cls}$	Classification Loss
$L_{usability}$	LLM’s advice usability loss
$L_{text}$	LLM text enhanced loss
$L_{distill}$	FACTGUARD-D’s distillation loss
$L_{total}$	FACTGUARD’s total loss
$L_{FGD}$	FACTGUARD-D’s total loss
$;$	Concat corresponding dimensions of vectors
$*$	Multiply vectors

Table 3: Notation and their corresponding descriptions.

## C Algorithm

The algorithms are divided into two parts. The first part, shown in Algorithm 1 and Algorithm 3, presents the FACTGUARD and FACTGUARD-D training procedures. The sec-

---

**Algorithm 1: FACTGUARD Training Process**

---

**Input:** News  $n$ , LLM extracted topic-content  $c$ , LLM commonsense rationale  $r$ , ground truth labels  $y$ , LLM commonsense rationale judgment  $y_{llm}$ , number of samples  $N$ , batch size  $B$

- 1: **Initialization:** hyperparameters  $\alpha, \beta$ , epoch  $E \leftarrow N/B$ , optimizer  $\leftarrow AdamW$
- 2: **for**  $epoch = 1$  to  $E$  **do**
- 3:    $C \leftarrow$  News Encoder( $n$ )
- 4:    $C, R \leftarrow$  LLM Text Encoder( $c, r$ )
- 5:    $\hat{y}_2 \leftarrow$  Classifier( $C$ )
- 6:    $\hat{y}_3 \leftarrow$  Classifier( $R$ )
- 7:   **for**  $i = 1$  to  $2$  **do**
- 8:      $f_{C \rightarrow Ri} \leftarrow CA(C, R, R)$
- 9:      $f_{R \rightarrow Ci} \leftarrow CA(R, C, C)$
- 10:      $w_i \leftarrow$  weight\_mapper( $f_{R \rightarrow Ci}$ )
- 11:      $\hat{w}_i \leftarrow$  MLP( $f_{R \rightarrow Ci}$ )
- 12:   **end for**
- 13:    $f_{llm} \leftarrow [w_1 * f_{C \rightarrow R1}; w_2 * f_{C \rightarrow R2}]$
- 14:    $f_N \leftarrow [Attn(N) + Attn(N)]/2$
- 15:    $f_{cls} \leftarrow [f_{llm}; f_N]$
- 16:    $\hat{y} \leftarrow$  MLP( $f_{cls}$ )
- 17:    $L_{cls} \leftarrow BCE(\hat{y}, y)$
- 18:    $L_{weight} \leftarrow BCE(\hat{w}_1, 0) + BCE(\hat{w}_2, y_{llm})$
- 19:    $L_{text} \leftarrow CE(\hat{y}_2, y) + CE(\hat{y}_3, y_{llm})$
- 20:    $L_{total} \leftarrow L_{cls} + \alpha \frac{L_{usability}}{2} + \beta \frac{L_{text}}{2}$
- 21:   **Zero gradients:**  $optimizer.zero\_grad()$
- 22:   **Backward pass:**  $L_{total}.backward()$
- 23:   **Update weights:**  $optimizer.step()$
- 24: **end for**

---

ond part, detailed in Algorithm 2 and Algorithm 4, illustrates the FACTGUARD and FACTGUARD-D inference processes.

## D Data Process

In this paper, we leverage an LLM for data processing through prompt engineering. The process consists of two main components: (1) extracting topics and core content to mitigate the influence of text style, and (2) generating commonsense rationales and judgments to infuse additional knowledge from the LLM. The overall workflow is illustrated in Figure 2.

### D.1 News Topic-Content Extraction

Given that the topic of a news article best encapsulates the overall meaning, and the main content centers on the core event—thus minimizing stylistic influence—our method employs LLMs to extract both the topic and main content of news texts. The detailed extraction process is illustrated in the left half of Figure 2. Prompt semantics are kept consistent across Chinese and English, with translations performed as needed. During extraction, we assess both the semantic similarity and information density between the extracted text and the original news to evaluate relevance and mitigate the risk of LLM hallucinations. Detailed experimental results are presented in Section D.2. Owing to

---

**Algorithm 2: FACTGUARD Inference Process**

---

**Input:** News  $n$ , LLM extracted topic-content  $c$ , LLM commonsense rationale  $r$

**Output:** Prediction  $\hat{y}$

- 1: **Initialization:** Frozen FACTGUARD model
- 2:  $C \leftarrow$  News Encoder( $n$ )
- 3:  $C, R \leftarrow$  LLM Text Encoder( $c, r$ )
- 4: **for**  $i = 1$  to  $2$  **do**
- 5:    $f_{C \rightarrow Ri} \leftarrow CA(C, R, R)$
- 6:    $f_{R \rightarrow Ci} \leftarrow CA(R, C, C)$
- 7:    $w_i \leftarrow$  MLP( $f_{R \rightarrow Ci}$ )
- 8: **end for**
- 9:  $f_{llm} \leftarrow [w_1 * f_{C \rightarrow R1}; w_2 * f_{C \rightarrow R2}]$
- 10:  $f_N \leftarrow [Attn(n) + Attn(n)]/2$
- 11:  $f_{cls} \leftarrow [f_{llm}; f_N]$
- 12:  $\hat{y} \leftarrow$  MLP( $f_{cls}$ )
- 13: **return**  $\hat{y}$

---

---

**Algorithm 3: FACTGUARD-D Training Process**

---

**Input:** News  $n$ , LLM extracted topic-content  $c$ , LLM commonsense rationale  $r$ , ground truth labels  $y$ , number of samples  $N$ , batch size  $B$

- 1: **Initialization:** Frozen FACTGUARD model, frozen News Encoder from FACTGUARD, frozen MLP from FACTGUARD, hyperparameter  $\lambda$ , epoch  $E \leftarrow N/B$ , optimizer  $\leftarrow AdamW$
- 2: **for**  $epoch = 1$  to  $E$  **do**
- 3:    $f_{cls} \leftarrow$  FACTGUARD( $c, r, n$ )
- 4:    $C \leftarrow$  News Encoder( $n$ )
- 5:    $f_{cls}^d \leftarrow$  Feature Simulator( $C$ )
- 6:    $\hat{y} \leftarrow$  MLP( $f_{cls}^d$ )
- 7:    $L_{cls} \leftarrow BCE(\hat{y}, y)$
- 8:    $L_{distill} \leftarrow MSE(f_{cls}, f_{cls}^d)$
- 9:    $L_{FGD} \leftarrow L_{cls} + \lambda L_{distill}$
- 10:   **Zero gradients:**  $optimizer.zero\_grad()$
- 11:   **Backward pass:**  $L_{FGD}.backward()$
- 12:   **Update weights:**  $optimizer.step()$
- 13: **end for**

---

differences in grammatical structure and information density between Chinese and English, the similarity scores for Weibo21 (Chinese) are typically lower than those for GossipCop (English). Accordingly, we require that the cosine similarity between extracted English news topic-content and the original news exceeds 0.9, while for extracted Chinese news topic-content the threshold is set at 0.8 according to the distribution of most similarity scores. If these criteria are not met, the LLM is prompted to regenerate the extraction use advanced LLM, such as Chatgpt-4o (OpenAI 2024) and DeepSeek-R1 (DeepSeek-AI 2025). In addition, after extraction, we further evaluate the information content of the results using Shannon entropy.

### D.2 Large Model Extraction Metrics

By evaluating both cosine similarity and Shannon entropy, we ensure that the extracted content is both faithful to the

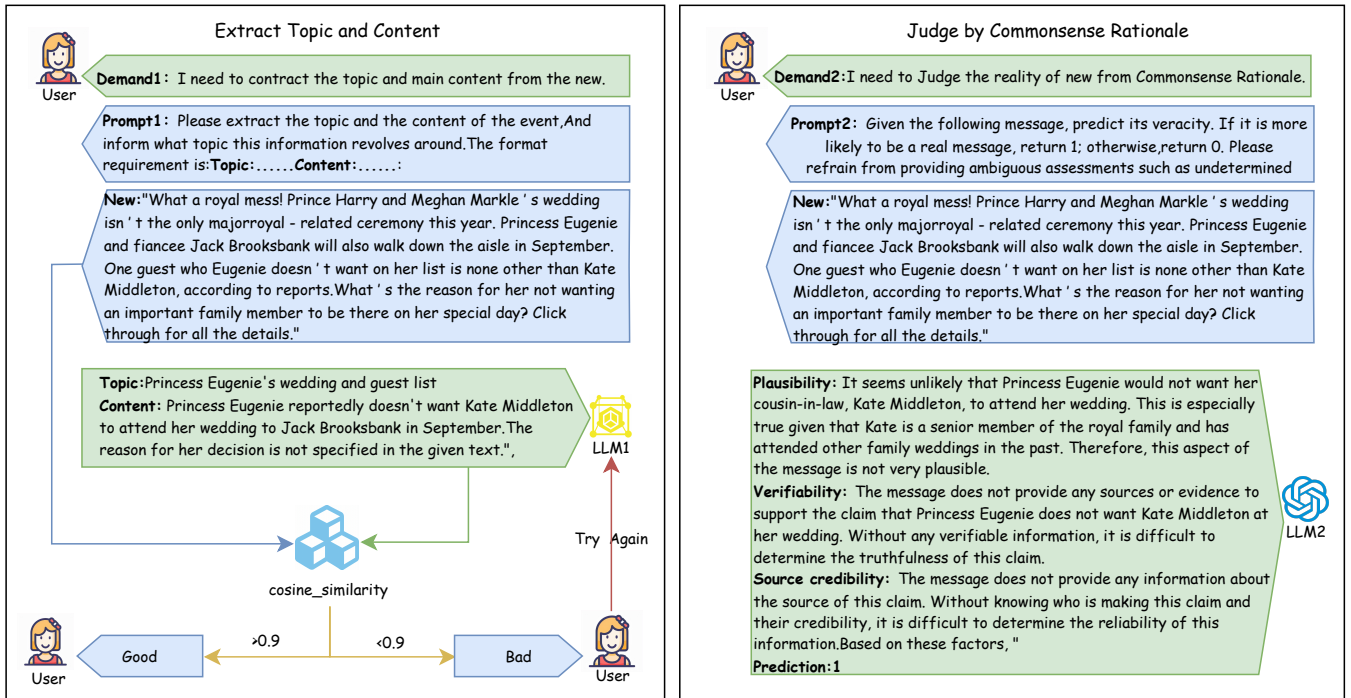


Figure 2: Data processing workflow. It consists of two parts: (1) extracting topics and content to mitigate the influence of text style; and (2) performing commonsense reasoning to identify contradictions in the news and generate LLM judgments.

#### Algorithm 4: FACTGUARD-D Inference Process

**Input:** News  $n$   
**Output:** Prediction  $\hat{y}$

- 1: **Initialization:** Frozen FACTGUARD-D model
- 2:  $C \leftarrow \text{News Encoder}(n)$
- 3:  $f_{cls}^d \leftarrow \text{Feature Simulator}(C)$
- 4:  $\hat{y} \leftarrow \text{MLP}(f_{cls}^d)$
- 5: **return**  $\hat{y}$

original news and sufficiently informative.

First, we input both the generated text and the real text into the corresponding Text Encoder, and assess their semantic similarity using cosine similarity. This approach helps constrain potential hallucinations from the LLM during content extraction. As shown in the top panels of Figure 3, the cosine similarity between the extracted topic-content, and the original news serves as a constraint during the extraction process—exceeding 0.8 for the Weibo21 dataset and 0.9 for the GossipCop dataset. If the similarity does not reach the specified threshold, the extraction process is repeated until the requirement is satisfied. Finally, across all data splits, the average cosine similarity between the extracted topic-content and the original news consistently exceeds 0.9 in two datasets. These results indicate that the extracted topics and content are highly semantically aligned with the original news, demonstrating that the LLM fulfills our requirements with minimal hallucination.

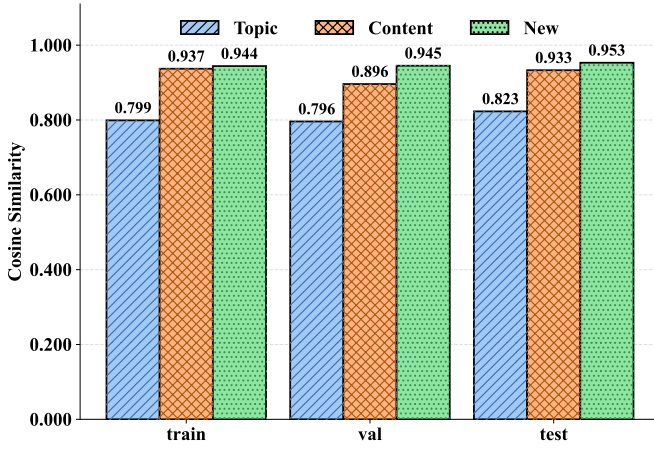
Next, we calculate the Shannon entropy for both the gen-

erated and real texts to evaluate whether the generated content retains sufficient information (Cao, Nguyen, and Zafarani 2025). A significantly lower Shannon entropy in the generated text compared to the real text may indicate overly simplified content and possible hallucination. The bottom panels of Figure 3 present the average entropy values per information piece in the datasets. Here, notable differences emerge between Weibo21 and GossipCop datasets. The original English news is longer and has a higher information density, yet after extraction, the information density drops to about half that of the original text. In contrast, the original Chinese news is shorter, and due to the inherently high information density of Chinese, the extracted topic-content maintains a similar information density as the original. Overall, the high semantic similarity and reasonable information density of the extracted content indicate that the LLM-based extraction process is both precise and reliable, with low risk of content hallucination across both Weibo21 and GossipCop datasets.

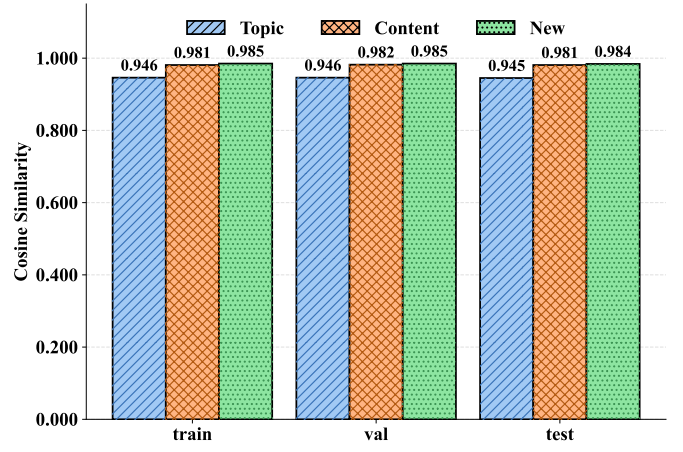
Through these two indicators, we screened and evaluated the content output by LLMs, and finally reduced the illusion of LLMs in the task of extracting news topic-content as much as possible.

### D.3 Commonsense Reasoning and Judgment

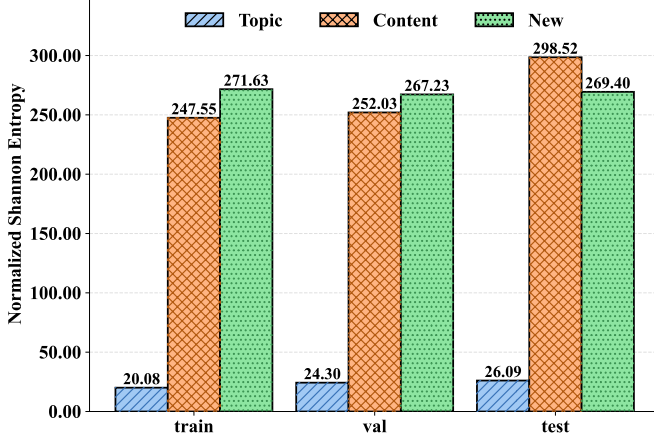
Commonsense reasoning leverages the inherent knowledge base of LLM to identify contradictions within news content from a commonsense perspective, aligning more closely with general human cognition. This component follows the approach proposed in ARG (Hu et al. 2024), with both the



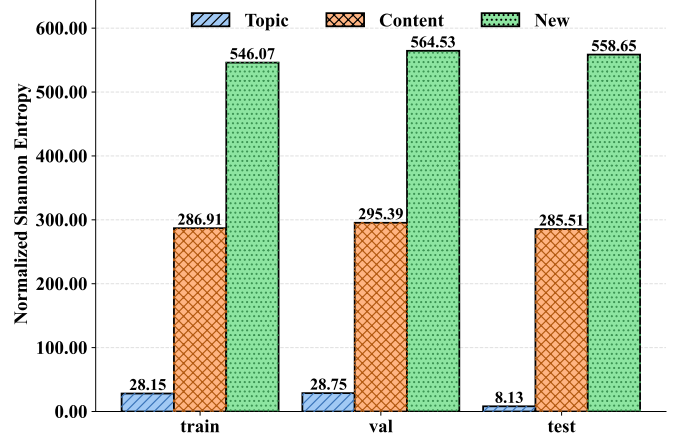
(a) Weibo21 dataset similarity



(b) GossipCop dataset similarity



(c) Weibo21 dataset Shannon Entropy



(d) GossipCop dataset Shannon Entropy

Figure 3: Similarity and Shannon entropy analysis on Weibo21 and GossipCop datasets.

Weibo21 and GossipCop datasets’ commonsense reasoning and judgment modules implemented using ChatGPT3.5-turbo. Prompt semantics are consistent across both languages, with translations as needed. The corresponding workflow is illustrated in the right half of Figure 2.

## E Experimental Setup

### E.1 Dataset

We employ the Weibo21 (Chinese) (Nan et al. 2021) and GossipCop (English) (Shu et al. 2020) for evaluation. Both datasets are preprocessed by deduplication and temporal splitting, following established practices (Zhu et al. 2022; Mu, Bontcheva, and Aletras 2023; Hu et al. 2024), to mitigate the risk of data leakage and prevent overestimation of SLM performance. In addition, we also utilize the commonsense rationales from (Hu et al. 2024). The statistical summaries of these datasets are provided in Table 4.

### E.2 Baselines

Recent advances in early fake news detection predominantly leverage LLMs and SLMs which only uses news content. In

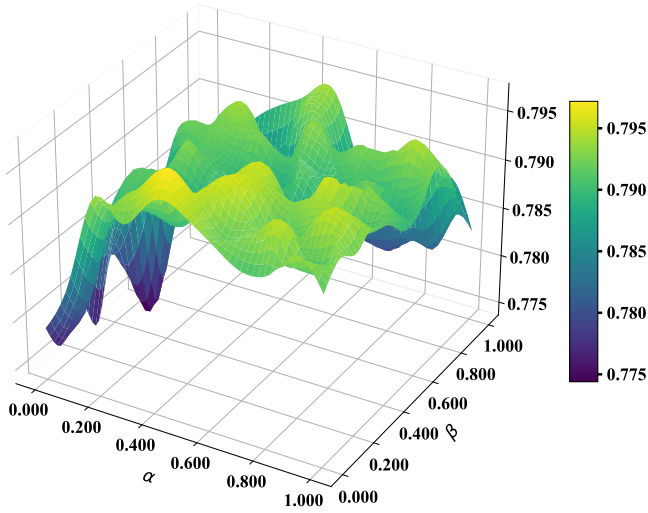
#	Weibo21			GossipCop		
	Train	Val	Test	Train	Val	Test
Real	2,331	1,172	1,137	2,878	1,030	1,024
Fake	2,873	779	814	1,006	244	234
Total	5,204	1,951	1,951	3,884	1,274	1,258

Table 4: Statistics of the number of real and fake samples in Weibo21 and GossipCop datasets.

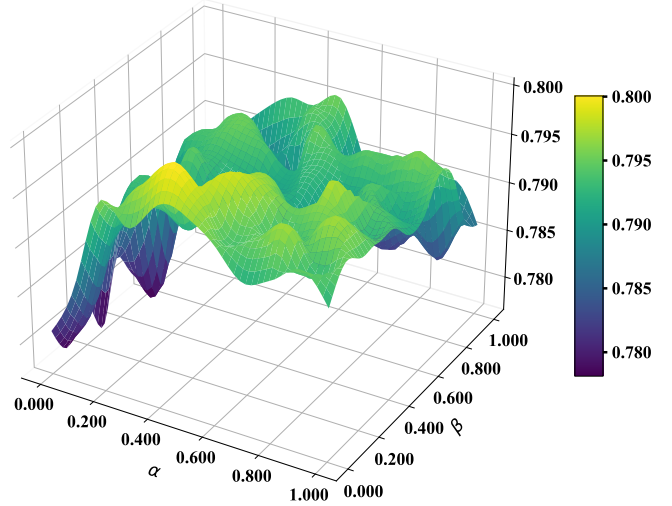
this study, we select 16 representative methods and categorize them into four groups: Group 1 comprises LLM-only methods; Group 2 consists of SLM-only methods; Group 3 includes LLM-SLM methods; and Group 4 focuses on knowledge distillation.

**G1: LLM-only.** These methods directly employ prompt engineering and multi-agent frameworks with LLMs for fake news detection.

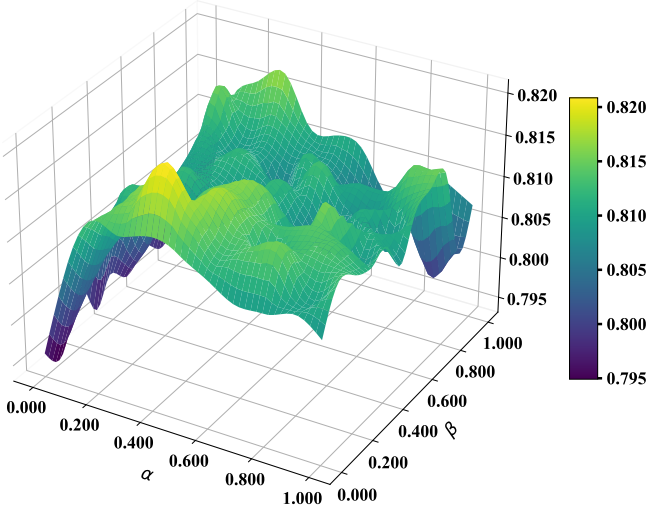
1. **GPT-3.5-turbo** (OpenAI 2023): Employed in conjunction with few-shot learning for Weibo21 dataset and few-shot CoT for GossipCop dataset.



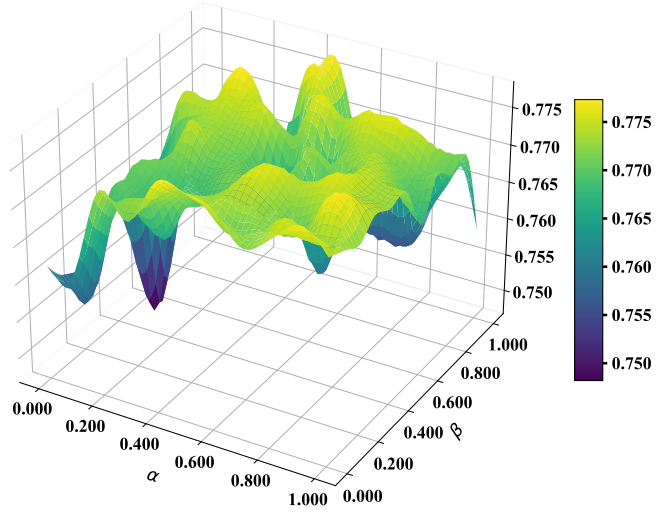
(a) the macF1 of Weibo21 dataset



(b) the Acc. of Weibo21 dataset



(c) the  $F1_{\text{real}}$  of Weibo21 dataset



(d) the  $F1_{\text{fake}}$  of Weibo21 dataset

2. **GPT-4o-mini** (OpenAI 2024): Direct use of GPT-4o-mini as an LLM detector for fake news detection.
3. **ChatEval-o** (Chan et al. 2024): Utilizes the one-by-one strategy for fake news migration within the multi-agent debate framework ChatEval.
4. **ChatEval-s** (Chan et al. 2024): Applies the Simultaneous-Talk strategy in the ChatEval multi-agent debate framework to tackle the fake news migration task.

**G2: SLM-only.** SLMs, such as BERT or RoBERTa, generally perform well on fake news detection tasks. This category includes:

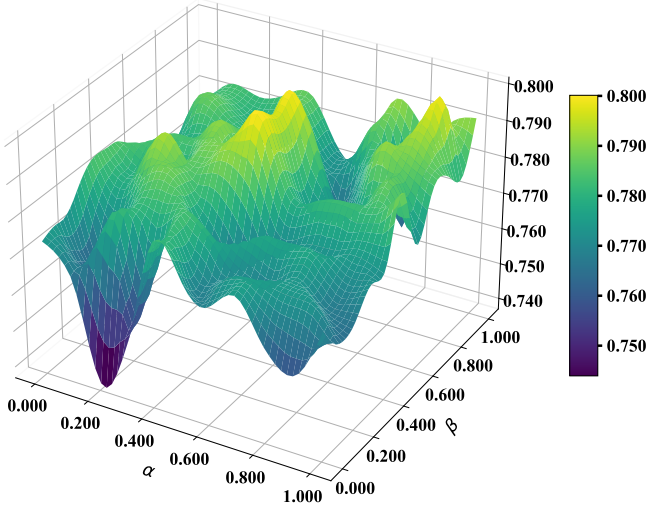
1. **BERT** (Devlin et al. 2019): Fake news detection using a fine-tuned vanilla BERT-base model.
2. **RoBERTa** (Liu et al. 2019): Employs RoBERTa as the text encoder, incorporating linear attention for enhanced fake news detection.

3. **EANN** (Wang et al. 2018): Employs auxiliary adversarial training to isolate and reduce event-related features, using the publication year as an auxiliary label.
4. **Publisher-Emo** (Zhang et al. 2021): Integrates emotional attributes with textual features for fake news detection.
5. **ENDEFA** (Zhu et al. 2022): Utilizes causal learning to eliminate entity bias, enhancing generalization under distribution shifts.

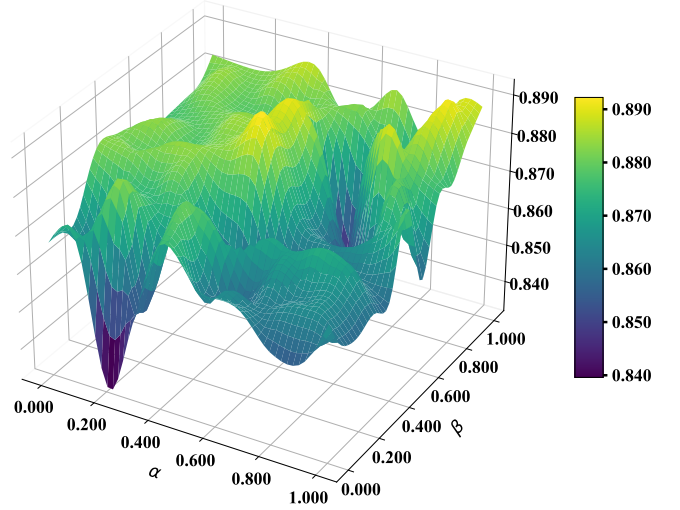
**G3: LLM-SLM.** LLMs alone often do not perform optimally. Therefore, several approaches combine LLM and SLM to enhance performance.

1. **BERT+Rationale** (Hu et al. 2024): Combines features from both the news and rationale encoders, feeding them into an MLP for prediction.
2. **SuperICL** (Zhong et al. 2023): Uses an SLM plug-in to enhance the in-context learning capabilities of the LLM

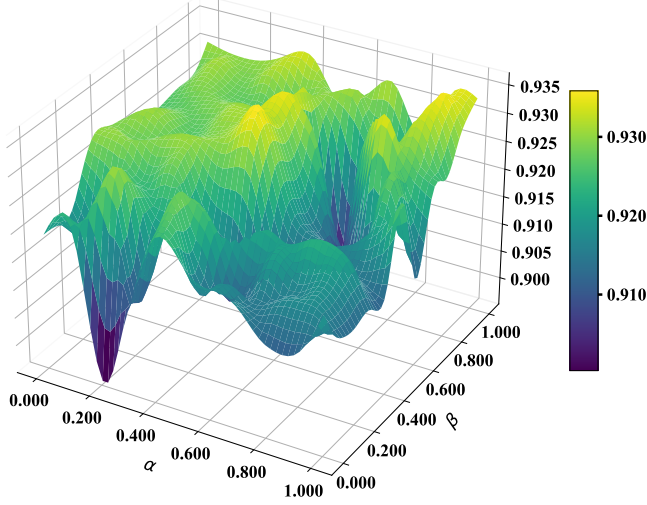




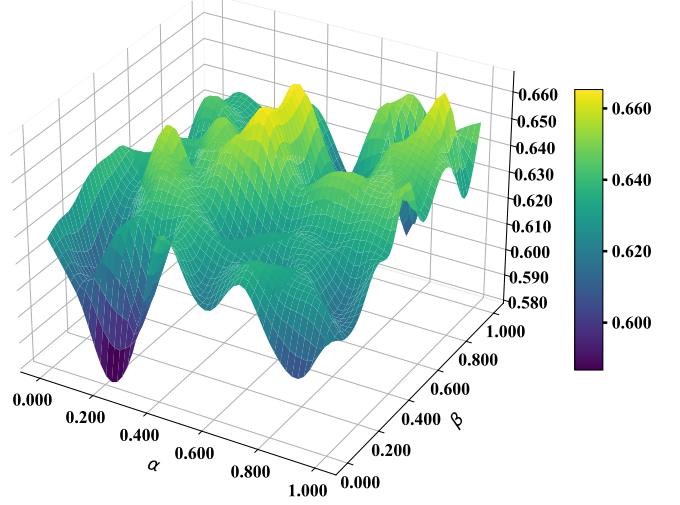
(e) the macF1 of GossipCop dataset



(f) the Acc. of GossipCop dataset



(g) the  $F1_{\text{real}}$  of GossipCop dataset



(h) the  $F1_{\text{fake}}$  of GossipCop dataset

Figure 4: Sensitivity analysis of FACTGUARD model in Weibo21 and GossipCop datasets across four evaluation metrics: macF1, Accuracy,  $F1_{\text{real}}$ , and  $F1_{\text{fake}}$ .

by incorporating predictions and confidence levels for each test sample into the prompt.

3. **BERT+GenFEND** (Nan et al. 2024): Combines features from both the news and LLMs’ comments based on BERT and qwen3-235b-a22b-instruct-2507<sup>8</sup> (Team 2025).
4. **RoBERTa+GenFEND** (Nan et al. 2024): Combines features from both the news and LLMs’ comments based on Roberta and qwen3-235b-a22b-instruct-2507.
5. **TED** (Liu et al. 2025): A multi-agent system leveraging LLM-driven structured debates to enhance both interpretability and accuracy.

<sup>8</sup><https://huggingface.co/Qwen/Qwen3-235B-A22B-Instruct-2507>

**G4: Distilled Model.** Distilled models are particularly suitable for resource-constrained and cold-start scenarios.

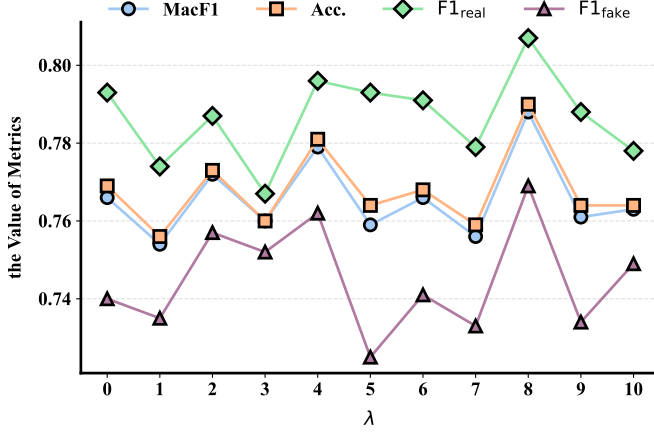
1. **ARG-D** (Hu et al. 2024): A rationale-free version of ARG created via distillation, designed for cost-sensitive applications that do not require LLM queries.

### E.3 Metrics

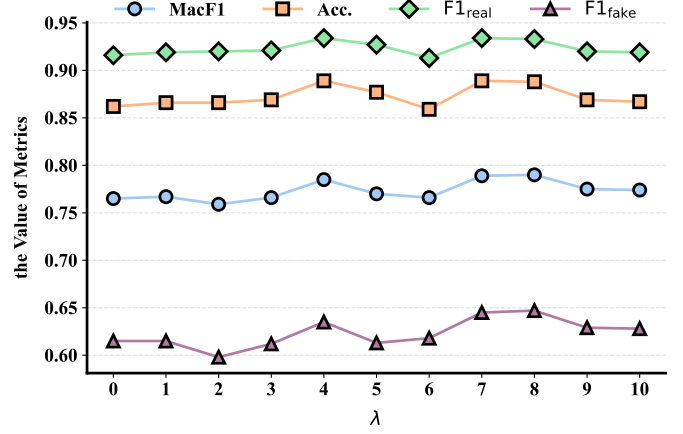
This study employs four evaluation metrics: Accuracy (Acc.),  $F1_{\text{real}}$ ,  $F1_{\text{fake}}$ , and Macro-F1 (macF1), to comprehensively assess the performance of the FACTGUARD model.

**Acc.** measures the proportion of all news samples (both real and fake) that are correctly classified by the model:

$$\text{Acc.} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}, \quad (20)$$



(a) the value of Weibo21 dataset distillation result



(b) the value of GossipCop dataset distillation result

Figure 5: Sensitivity analysis of FACTGUARD-D model in Weibo21 and GossipCop datasets across four evaluation metrics: macF1, Acc., F1<sub>real</sub>, and F1<sub>fake</sub>.

where TP (True Positive) is the number of fake news articles correctly identified, TN (True Negative) is the number of real news articles correctly identified, FP (False Positive) is the number of real news articles incorrectly classified as fake, and FN (False Negative) is the number of fake news articles incorrectly classified as real. Together, these four quantities constitute the confusion matrix, providing a comprehensive view of the model’s classification performance.

**F1<sub>real</sub>** is the F1 score for the “real news” category, reflecting the model’s performance in identifying real news:

$$F1_{real} = 2 \times \frac{\text{Precision}_{real} \times \text{Recall}_{real}}{\text{Precision}_{real} + \text{Recall}_{real}}, \quad (21)$$

where  $\text{Precision}_{real}$  denotes the proportion of news predicted as real that is actually real, and  $\text{Recall}_{real}$  represents the proportion of true real news that is correctly identified.

**F1<sub>fake</sub>** is the F1 score for the “fake news” category, emphasizing the model’s detection ability for fake news:

$$F1_{fake} = 2 \times \frac{\text{Precision}_{fake} \times \text{Recall}_{fake}}{\text{Precision}_{fake} + \text{Recall}_{fake}}, \quad (22)$$

where  $\text{Precision}_{fake}$  denotes the proportion of news predicted as fake that is actually fake, and  $\text{Recall}_{fake}$  indicates the proportion of true fake news that is successfully detected.

**macF1** is the average F1 score across both categories, providing an overall measure of recognition capability:

$$\text{macF1} = \frac{F1_{fake} + F1_{real}}{2}. \quad (23)$$

## F Sensitivity Study

This section presents parameter sensitivity experiments for two loss function hyperparameters as well as the choice of text encoder of FACTGUARD and FACTGUARD-D to identify the optimal model configuration.

### F.1 Loss Hyperparameters

The loss function of FACTGUARD involves two key hyperparameters,  $\alpha$  and  $\beta$ . We employ a grid search to determine their optimal values. As illustrated in Figure 4, we evaluate four performance metrics (i.e. Acc., F1<sub>real</sub>, F1<sub>fake</sub> and macF1) on both Weibo21 and GossipCop datasets under various hyperparameter settings. In the plots, the x-axis denotes  $\alpha$ , the y-axis denotes  $\beta$ , and the z-axis represents the value of the corresponding metric. Both  $\alpha$  and  $\beta$  are discretized into 11 grid points from 0 to 10. Figure 4 displays the results of this grid search.

The results indicate that FACTGUARD achieves optimal performance on the Weibo21 dataset when  $\alpha = 0.4$  and  $\beta$  is within the range (0.1, 0.2). For the GossipCop dataset, the best results of FACTGUARD are observed when  $\alpha = 0.5$  and  $\beta$  falls within the range (0.5, 0.6). Based on these findings, we further conduct a random search for  $\beta$  and ultimately determine the optimal hyperparameters of FACTGUARD to be  $\alpha = 0.40$ ,  $\beta = 0.16$  in Chinese FACTGUARD, and  $\alpha = 0.50$ ,  $\beta = 0.58$  in English FACTGUARD.

The loss function of FACTGUARD-D involves a key hyperparameter, the distillation coefficient  $\lambda$ . As shown in Figure 5, we evaluate eight performance metrics on both Weibo21 and GossipCop datasets under different hyperparameter settings. In the plots, the x-axis represents  $\lambda$ , while the y-axis indicates the value of the corresponding metric. The value of  $\lambda$  is discretized into 11 points ranging from 0 to 10. Figure 5 presents the results of this sequential hyperparameter search. The results demonstrate that the optimal performance on both the Weibo21 and GossipCop datasets is achieved when  $\lambda = 8$ .

### F.2 Text Encoder

Due to differences in language characteristics between Chinese and English, the choice of text encoder can significantly affect FACTGUARD performance. The comparative results for BERT and RoBERTa are summarized in Ta-

Text Encoder	Weibo21				GossipCop			
	macF1	Acc.	F1 <sub>real</sub>	F1 <sub>fake</sub>	macF1	Acc.	F1 <sub>real</sub>	F1 <sub>fake</sub>
Bert	<b>0.801</b>	<b>0.804</b>	<b>0.824</b>	<b>0.777</b>	0.784	0.879	0.927	0.642
RoBERTa	0.764	0.765	0.784	0.745	<b>0.805</b>	<b>0.892</b>	<b>0.935</b>	<b>0.675</b>

Table 5: the performance of BERT and RoBERTa as FACTGUARD’s encoder.

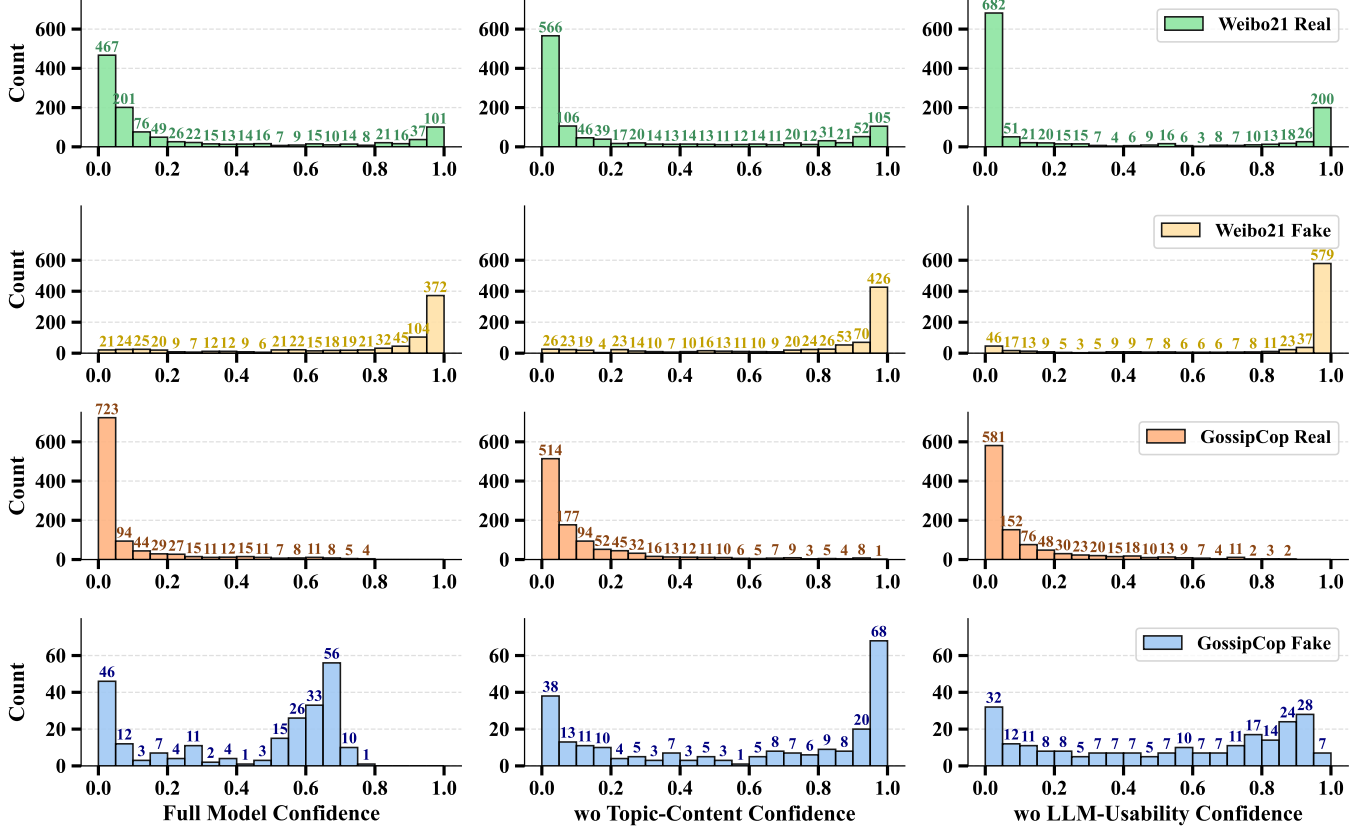


Figure 6: Confidence distributions of real and fake samples in GossipCop and Weibo21 dataset before and after with style debiasing and llm usability judgment in FACTGUARD.

ble 5. For Chinese FACTGUARD, BERT is sourced from Google’s bert-base-chinese<sup>9</sup> (Devlin et al. 2019) pretrained model, and RoBERTa from hfl’s chinese-roberta-wwm-ext<sup>10</sup> (Cui et al. 2020) pretrained model. For English FACTGUARD, BERT is based on Google’s bert-base-uncased pretrained model<sup>11</sup>, and RoBERTa from Facebook’s roberta-base<sup>12</sup> (Liu et al. 2019) pretrained model. Based on the experimental results, Google’s BERT is selected as the Chinese FACTGUARD text encoder and Facebook’s RoBERTa is chosen for English FACTGUARD.

### E.3 Confidence Distribution Analysis

Figure 6 depicts the confidence distributions of real and fake samples on the Weibo21 and GossipCop datasets under three model configurations: the full FACTGUARD model, the ablation without topic-content features, and the ablation without LLM usability evaluator. This comparison allows us to disentangle the unique contributions of style debiasing and LLM-based factual judgment to model calibration and detection performance.

Across both datasets, we observe that the full FACTGUARD model, which integrates both style debiasing (via topic-content disentanglement) and LLM usability evaluation, yields the most desirable and interpretable confidence distributions. For fake news samples, the full FACTGUARD model markedly suppresses the occurrence of extremely high-confidence predictions (i.e., fake samples prediction probabilities from 0.8 to 1). This moderation of confidence

<sup>9</sup><https://huggingface.co/google-bert/bert-base-chinese>

<sup>10</sup><https://huggingface.co/hfl/chinese-roberta-wwm-ext>

<sup>11</sup><https://huggingface.co/google-bert/bert-base-uncased>

<sup>12</sup><https://huggingface.co/FacebookAI/roberta-base>



**Case 1: Although the LLM made an incorrect prediction, FACTGUARD successfully detect fake news due to the influence of the style-debiased topic-content extracted by the LLM in Weibo21.**

**News:** #High Risk Of COVID Spread In Beijing# The COVID-19 situation in Beijing is severe. Academician Zhong Nanshan, aged 84, has once again gone to Beijing. He is the brightest star—salute to him! #Beijing Demands Tracing All Close Contacts In Shortest Time# — Video by Weibo user

**LLM Commonsense Rationale:** Based on the current COVID-19 situation, there is indeed a risk of virus spread in Beijing, and there are reports demanding that all close contacts be traced. However, it cannot be confirmed whether Academician Zhong Nanshan has once again gone to Beijing, nor can the content of the Weibo video be verified.

**LLM extracted topic-content:** Topic: Academician Zhong Nanshan’s Involvement in Beijing’s COVID-19 Prevention and Control Efforts. Content: Academician Zhong Nanshan has once again gone to Beijing. He is the brightest star—salute to him! The COVID-19 situation in Beijing is severe, with a high risk of spread. Authorities demand tracing all close contacts in the shortest possible time. This is shown in a video by Weibo user.

**Ground Truth:** Fake    **LLM Judgment:** Real    **Full FACTGUARD:** Fake    **FACTGUARD w/o Topic-Content:** Real  
**Full FACTGUARD Confidence:** 0.808    **FACTGUARD w/o Topic-Content Confidence:** 0.459

**Case 2: Although the LLM made an correct prediction, FACTGUARD without the usability module fails to effectively detect fake news due to its inability to integrate the recommendations from LLM in Weibo21.**

**News:** In the video, it is actor Zhang Fengyi’s grandson. He was bitten by his own dog and thought it was nothing, so he didn’t go to the hospital for a rabies vaccine. In the end, he got rabies and barked like a dog. People who keep dogs and cats must pay special attention. Don’t be harmed by the very pet parents you value!

**LLM Commonsense Rationale:** Plausibility: The message seems plausible as it talks about a celebrity couple and their relationship issues, which is a common topic in the media. Verifiability: The message mentions exclusive details, but it is not clear who the source of this information is. Without a credible source, it is difficult to verify the accuracy of the information. Source credibility: The message does not provide any information about the source of the information, which makes it difficult to assess the credibility of the message.

**LLM extracted topic-content:** Topic: The Importance of Vaccinations for Cats and Dogs. Content: The video tells the story of actor Zhang Fengyi’s grandson, who contracted rabies after being bitten by the family dog and failing to receive a timely rabies vaccination. This incident serves as a reminder for pet owners to take the health management of their cats and dogs seriously, ensuring timely vaccinations to prevent avoidable health risks caused by negligence.

**Ground Truth:** Fake    **LLM Judgment:** Fake    **Full FACTGUARD:** Fake    **FACTGUARD w/o Topic-Content:** Real  
**Full FACTGUARD Confidence:** 0.740    **FACTGUARD w/o Topic-Content Confidence:** 0.113

is not merely a numerical adjustment; rather, it reflects the model’s enhanced ability to avoid overfitting to superficial stylistic artifacts and spurious correlations that are often present in fake news. Importantly, this reduction in overconfident judgments does not compromise overall detection performance; on the contrary, the model demonstrates improved effectiveness in identifying fake news, indicating that the learned decision boundary is more robust and generalizable. The resulting confidence scores for fake news become more evenly distributed, suggesting better calibration and a lower risk of making extreme, potentially erroneous decisions.

For fake news samples, when we ablate the topic-content module, thereby removing the style debiasing mechanism, a clear shift in the confidence distribution emerges. The model tends to revert to overconfident predictions for fake samples, with a noticeable increase in the number of samples assigned near-certain fake probabilities. This pattern indicates a renewed reliance on shallow stylistic or topical cues that are frequently entangled with fake news, leading to poorer generalization. The absence of topic-content disentanglement thus exposes the model to the risk of exploiting dataset-specific artifacts, underlining the critical role of style debiasing in mitigating such bias and promoting more reliable,

fact-oriented detection.

Ablating the LLM usability component yields a similarly instructive outcome. Without the guidance of LLM-based factual judgment, the model again displays a tendency toward extreme confidence for fake news and less decisive recognition of real news. The removal of this module results in a less nuanced decision process: fake samples are concentrated at the high-confidence end, while real samples suffer from a reduction in high-confidence, correct predictions and a concomitant increase in low-confidence assignments. This highlights the importance of LLM usability in enforcing a fact-centric evaluation framework, which complements style debiasing by promoting decisions grounded in semantic and factual consistency rather than surface-level features.

For real news samples, these trends are mirrored but in the opposite direction. The full FACTGUARD model is able to assign a higher proportion of real samples with strong confidence for being real(i.e., real samples prediction probabilities from 0 to 0.2), while the number of low-confidence assignments decreases. This shift reflects the model’s increased decisiveness and accuracy, which is diminished when either the style debiasing or LLM usability module is removed. Notably, the topic-content ablation leads to greater uncertainty and misclassification for real news, reaffirming

<b>Case 3: Although the LLM made an incorrect prediction, FACTGUARD successfully detect fake news due to the influence of the style-debiased topic-content extracted by the LLM in GossipCop.</b>			
<p><b>News:</b> Is it even a proper Independence Day celebration without fireworks, Stars and Stripes face paint, and photos from one of Taylor Swift’s infamous Fourth of July parties? Yes it is, but still the Internet had to do without the latter in 2018, because this was the second year in a row that her home in Watch Hill, Rhode Island was empty on America’s birthday. ( In other words, perhaps Taylor not celebrating the 4th during Trump’s presidency is her biggest political statement yet? ) Advertisement - Continue Reading Below</p> <p>Instead of holding one of her # squad - powered summer bashes this year, T - Swift chose to spend a break from her current world tour in Turks and Caicos with her boyfriend actor Joe Alwyn. The couple were spotted all loved up and canoodling on the beach.</p> <p><b>LLM Commonsense Rationale:</b> Plausibility: It is plausible that Taylor Swift would host a Fourth of July party, as she has done so in the past. It is also plausible that she would choose to spend the holiday with her boyfriend instead. Verifiability: The message provides specific details about Taylor Swift’s whereabouts on the Fourth of July, including her absence from her home in Watch Hill, Rhode Island and her vacation in Turks and Caicos with her boyfriend. These details could potentially be verified through other sources. Source credibility: The source of the message is not provided, so it is unclear how credible it is. Based on these factors.</p> <p><b>LLM extracted topic-content:</b> Topic: Taylor Swift’s absence from Fourth of July celebrations. Content: Taylor Swift did not host her annual Fourth of July party for the second year in a row. Instead, she spent time with her boyfriend Joe Alwyn in Turks and Caicos. This absence could be interpreted as a political statement against Trump’s presidency.</p>			
<b>Ground Truth:</b> Fake <b>LLM Judgment:</b> Real <b>Full FACTGUARD:</b> Fake <b>FACTGUARD w/o Topic-Content:</b> Real <b>Full FACTGUARD Confidence:</b> 0.569 <b>FACTGUARD w/o Topic-Content Confidence:</b> 0.195			
<b>Case 4: Although the LLM made an uncertain prediction, FACTGUARD successfully detect fake news due to the influence of the LLM rationale usability module in GossipCop.</b>			
<p><b>News:</b> Whew! Rihanna and Hassan Jameel finally proved they are still together. We ’ ve got EXCLUSIVE details on what ’ s behind their make ups and break ups after they were photographed fighting. Fans were relieved when pics of Rihanna surfaced on a sexy Mexican vacation with boyfriend of one year Hassan Jameel. The couple is still together despite breakup reports, but RiRi looked like she was reading the Saudi billionaire the riot act in the photos, with her hand animatedly raised in the and she appeared to be yelling at him. “ Rihanna and Hassan have an extremely passionate and intense relationship, with a tendency to break - up and make - up fairly frequently. They have a very strong physical connection, and when they fight it can get pretty heated, but then they have crazy.</p> <p><b>LLM Commonsense Rationale:</b> Plausibility: The message seems plausible as it talks about a celebrity couple and their relationship issues, which is a common topic in the media. Verifiability: The message mentions exclusive details, but it is not clear who the source of this information is. Without a credible source, it is difficult to verify the accuracy of the information. Source credibility: The message does not provide any information about the source of the information, which makes it difficult to assess the credibility of the message.</p> <p><b>LLM extracted topic-content:</b> Topic: Rihanna and Hassan Jameel’s Relationship. Content: Rihanna and her boyfriend Hassan Jameel were spotted on a Mexican vacation, proving they are still together despite breakup rumors. They have a passionate and intense relationship with frequent breakups and makeups. Their fights can get heated, but they have a strong physical connection.</p>			
<b>Ground Truth:</b> Fake <b>LLM Judgment:</b> Other <b>Full FACTGUARD:</b> Fake <b>FACTGUARD w/o Topic-Content:</b> Real <b>Full FACTGUARD Confidence:</b> 0.579 <b>FACTGUARD w/o LLM Usability Confidence:</b> 0.480			

Figure 7: Case analysis: style-debiased topic-content and llm usability modules boost FACTGUARD’s fake news detection.

the importance of style disentanglement for robust recognition of authentic information.

In summary, the experimental results across both Weibo21 and GossipCop datasets consistently demonstrate that the integration of topic-content style debiasing and LLM usability evaluation not only optimizes detection performance but also produces confidence distributions that are more interpretable and reliable. The topic-content module plays a pivotal role in mitigating stylistic and topical biases, thereby preventing overfitting and promoting generalization, while LLM usability injects essential factual discernment into the decision process. The full FACTGUARD model, benefiting from both mechanisms, achieves superior robustness, calibration, and real-world applicability in automated news verification.

## F.4 Case Analysis

Figure 7 presents a comprehensive analysis of the impact of style-debiased topic-content extraction and the usability module within LLMs in automated fake news detection. Across both GossipCop and Weibo21 datasets, these modules together significantly enhance the reliability and interpretability of detection results.

Cases 1 and 3 illustrate how style-debiased topic-content extraction helps distill the factual core of news reports, even when the LLM initial prediction is incorrect. By filtering out emotional or exaggerated language, this module enables FACTGUARD to focus on verifiable information and more accurately identify falsehoods. For example, in Case 1 (Weibo21), although the LLM incorrectly predicts the news as real, the topic-content extraction isolates the uncertainty about whether Zhong Nanshan actually traveled to Beijing,

reducing reliance on sensational narrative and leading FACTGUARD to the correct “fake” classification with much higher confidence. In Case 3 (GossipCop), even after the LLM mislabels the news as real, the extraction process highlights the factual absence of Taylor Swift’s Fourth of July party, helping FACTGUARD ignore the gossip and focus on checkable facts, thereby improving robustness and accuracy in fake news detection.

Cases 2 and 4 highlight the limitations of the model when the LLM usability module is absent, making systematic evaluation of news credibility and verifiability difficult. In Case 2 (Weibo21), although the LLM correctly predicts the news as fake, without the usability module, FACTGUARD fails to effectively detect key fake news signals, such as the lack of source credibility and unclear attribution, resulting in very low confidence. Similarly, in Case 4 (GossipCop), when the LLM judgment is uncertain (“Other”), the removal of the usability module leads to the failure in identifying the absence of authoritative sources or verifiable evidence in the story about Rihanna and Hassan Jameel, lowering both confidence and interpretability. In both cases, the usability module is essential for integrating rationales about source reliability and verifiability, substantially boosting FACTGUARD’s confidence and final judgment.

Overall, the synergy between topic-content extraction and usability analysis modules allows FACTGUARD to compensate for limitations in direct LLM predictions. By enhancing both accuracy and interpretability, this modular approach demonstrates clear practical value for large-scale, automatic fake news detection in real-world scenarios.