

Reading the Lines, Decoding the Minds: Explainable Satirical Cartoon Detection via On-Canvas and Beyond-Canvas Modeling

Anonymous ACL submission

Abstract

Multimodal Satire Detection (MSD) aims to identify implicit criticism and stance in social media. However, existing MSD benchmarks mainly use user generated content and predominantly consist of realistic photographs, where satire is often conveyed through overt image text incongruity. By contrast, increasingly popular editorial cartoons convey satire through symbolic and metaphorical cues, requiring models to infer latent intent and pragmatics. To address these limitations, we propose MSCE, a dedicated evaluation benchmark for Multimodal Satirical Cartoons Evaluation. As a strong baseline for MSCE, we introduce a novel explainable framework named DOUBLE. DOUBLE endows Multimodal Large Language Model (MLLM) with a Clue2View mechanism to explicitly model interpretation through two complementary views: a literal on-canvas view and a metaphorical beyond-canvas view. In addition, DOUBLE incorporates a lightweight Small Language Model Arbiter to distill reasoning traces from the MLLM, ensuring reliable predictions with lower computational costs. Experiments on MSCE demonstrate that DOUBLE achieves the best performance, showcasing its ability to provide clear, well-grounded rationales for complex satirical cartoons.

1 Introduction

Satire has become a prevalent rhetorical device on social media, enabling users to comment on public affairs and signal stance. Satirical posts often employ overtly incongruent expressions to implicitly criticize or attack a target (Farabi et al., 2024; Tiwari et al., 2023). Effective identification and interpretation of multimodal satire is therefore consequential for a range of downstream applications, including misinformation detection (Rubin et al., 2016; Bedard and Schoenthaler, 2018; De Sarkar et al., 2018), public opinion monitoring (Lu et al., 2025), and sentiment analysis (Joshi et al., 2017; Mishra et al., 2016).

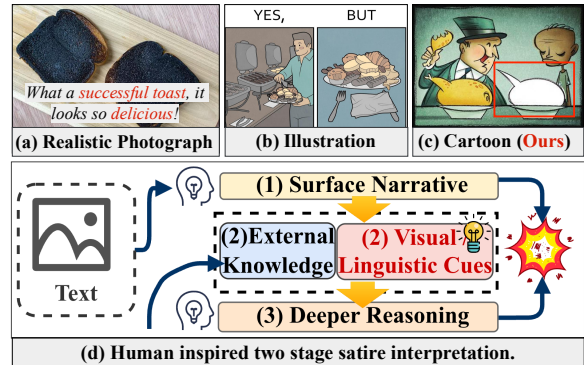


Figure 1: Comparison of image types in multimodal satire detection datasets. & Human inspired two stage satire interpretation.

Recently, the community has established several multimodal satire benchmarks on open social platforms, such as MMSD (Cai et al., 2019), MMSD2.0 (Qin et al., 2023), and SarcNet (Yue et al., 2024). They also have facilitated standardized evaluation and rapid methodological progress. A common assumption in prior work is that cross modal incongruity is a primary cue for satire (Pan et al., 2020; Liang et al., 2021). Models can leverage this incongruity to detect satire and achieve state of the art performance in open domain social media (Liang et al., 2022; Qiao et al., 2023; Wang et al., 2024a).

However, when the scope of investigation extends to satirical cartoons, existing benchmarks, which primarily feature realistic photographs or illustrations, exhibit notable limitations in two respects. (1) Most existing datasets are simply and directly collected from social media user generated content, where creators often elicit satire via explicit image text conflicts. As illustrated in Figure 1(a) and (b), the satire is easy to infer because it hinges on an explicit mismatch, either the caption praising “delicious” food despite burnt toast or the stark contrast between excessive stockpiling and subsequent waste. However, when satire is conveyed primarily through such obvious surface-level

071 cues, models may shortcut the reasoning process
072 and simply learn to detect salient anomalies. (2)
073 Recently, cartoons are a widely used for satire in
074 news commentary and public discourse (Kundu
075 et al., 2025). Crucially, cartoon satire rarely relies
076 on overt cues. Instead, it emerges from the com-
077 bination of symbolism and metaphor, suggestive
078 narratives, background knowledge about events and
079 eras, and intent inference. As a result, when satiri-
080 cal cues are veiled and the cartoon conveys an un-
081 derlying critical stance, models are more prone to
082 misclassification. For example, Figure 1(c) con-
083 veys satire through symbols (top hat & suit →
084 capitalist; frail body → working class), metaphor
085 (devouring → predatory exploitation), and context
086 (empty speech bubble → hollow promises). There-
087 fore, existing Multimodal Satire Detection (MSD)
088 benchmarks may fall short in evaluating non contra-
089 diction based satire, motivating a dedicated bench-
090 mark for editorial cartoons.

091 To address these limitations, we introduce
092 **MSCE**, the first Multimodal Satirical Cartoon
093 Evaluation benchmark, which contains 2,016 image-
094 text pairs collected from satirical cartoons
095 published in the *Satire and Humor* supplement
096 of *People’s Daily*. MSCE focuses on cartoons
097 whose image and text are surface consistent, where
098 satire arises from implicit intent rather than shallow
099 conflicts. Furthermore, it evaluates core abilities
100 for MSD, including symbol recognition, metaphor
101 mapping, contextual and era specific knowledge
102 use, and intent conflict inference. However, most
103 existing satire detection approaches treat cross
104 modal incongruity as the primary cue, and they
105 encounter two major challenges on MSCE as fol-
106 lows:

107 **C1: From surface consistency to metaphorical**
108 **reasoning.** Unlike conventional MSD, where satire
109 is typically signaled by explicit cross modal in-
110 congruity, editorial cartoons often maintain sur-
111 face consistency between text and image. And the
112 surface consistency conceals the underlying satiri-
113 cal intent. Detecting such satire therefore requires
114 metaphorical understanding beyond surface seman-
115 tics. In particular, models must perform multi-hop
116 reasoning that maps concrete visual symbols, such
117 as a broken bowl, to abstract concepts like eco-
118 nomic crisis, a capability that contradiction-based
119 approaches fail to capture.

120 **C2: Complexity of disentangling Literal Nar-**
121 **rative and Satirical Intent.** In editorial cartoons,
122 satirical cues are interwoven with literal narratives,

123 so indiscriminate fusion can let salient but non de-
124 cisive signals dominate and increase misclassifica-
125 tion. Even MLLMs may resort to coherent post
126 hoc stories without isolating true satirical triggers.
127 Therefore, effective satire detection requires ex-
128 plicit multi-perspective modeling: one encodes
129 the on-canvas literal narrative, another captures
130 beyond-canvas contextual intent, and the final pre-
131 diction is derived from their contrast.

132 To address the challenges, we propose
133 **DOUBLE**, a novel cartoon-oriented satire de-
134 tection framework powered by **Deriving On-**
135 **canvas Understanding and Beyond-canvas Logic**
136 **Extraction.** **DOUBLE** can explicitly models
137 the reasoning process to yield clear, evidence-
138 grounded rationales. It is inspired by the progres-
139 sive way humans interpret satirical cartoons (see
140 Figure 1 (d)). (1) Readers first form a surface nar-
141 rative from on-canvas content, and (2) then use
142 salient cues to retrieve beyond-canvas context to
143 (3) infer a deeper target and intent. Consequently,
144 satire arises from the contrast between the surface
145 reading and the context-enriched interpretation. Ac-
146 cordingly, we introduce **Clue2View**, a dual-view
147 mechanism that generates: (i) a literal *non-satirical*
148 explanation grounded in on-canvas cues, and (ii) a
149 context-enriched *satirical* explanation by integrat-
150 ing beyond-canvas knowledge, enabling models to
151 move beyond surface consistency and better detect
152 metaphor-driven satire.

153 Although fine-tuning multimodal large language
154 models (MLLMs) with **Clue2View** substantially
155 improves performance on MSD, it incurs expen-
156 sive computational cost and limits practical usabil-
157 ity. To address this issue, we further propose a
158 **Small Language Model Arbiter (SLM Arbiter)**
159 that serves as an efficient and reliable decision
160 maker within **DOUBLE**. Built upon a small lan-
161 guage model, **SLM Arbiter** distills knowledge from
162 **Clue2View** to refine the outputs of MLLMs through
163 lightweight fine-tuning, ensuring reliable predic-
164 tions with low computational cost. We summarize
165 our main contributions as follows.

- 166 • We further introduce **MSCE**, the first bench-
167 mark for satirical cartoon detection that fo-
168 cuses on veiled multimodal satire and requires
169 deep contextual and symbolic reasoning.
- 170 • We propose **DOUBLE**, an explainable frame-
171 work for MSD that is inspired by this progres-
172 sive human reasoning process. As a strong

173	baseline on MSCE, DOUBLE enables reliable	221
174	detection even when image and text are	222
175	surface consistent, while providing verifiable	223
176	evidence.	224
177	• Extensive experiments and case studies on	225
178	MSCE show that DOUBLE substantially out-	226
179	performs existing approaches for MSD, and	227
180	it establishes a strong baseline for future re-	228
181	search.	229
182	The code and data for our proposed DOUBLE	230
183	and MSCE are available at https://anonymous.	231
184	4open.science/r/DOUBLE-5577/ .	232
185	2 Related Work	233
186	2.1 Multimodal Satire Detection Benchmarks	234
187	Various benchmarks have been introduced for	235
188	MSD. Cai et al. (Cai et al., 2019) formulated	236
189	MSD on Twitter image–text posts and released	237
190	the MMSD dataset. Qin et al. (Qin et al., 2023)	238
191	subsequently refined MMSD by removing spuri-	239
192	ous cues and proposed MMSD2.0. Further, Desai	240
193	et al. (Desai et al., 2022) introduced the task of	241
194	multimodal satire explanation and constructed the	242
195	MORE dataset, where each instance is equipped	243
196	with natural language explanation. To evaluate the	244
197	capabilities of vision language models, the YesBut	245
198	benchmark adopted a multi task setup to assess	246
199	whether models truly understand and can explain	247
200	satire (Nandy et al., 2024). In the cross lingual	248
201	setting, SarcNet (Yue et al., 2024) provides English	249
202	and Chinese image–text pairs, complementing mul-	250
203	tilingual evaluation for MSD.	251
204	Most existing datasets rely primarily on realis-	252
205	tic photographic images, which tends to simplify	253
206	satire understanding. Even YesBut, despite includ-	254
207	ing cartoonist produced illustrations, offers limited	255
208	coverage of core cartoon satire mechanisms, includ-	256
209	ing symbolism, metaphor, and contextual inference.	257
210	To address these gaps, we introduce a benchmark	258
211	of multimodal satirical cartoons that emphasizes	259
212	implicit satire and demands complex contextual re-	260
213	asoning, aiming at promoting context driven satire	261
214	understanding.	262
215	2.2 Multimodal Satire Detection Methods	263
216	Most existing approaches to MSD are grounded in	264
217	modality incongruity, assuming that satire arises	265
218	from semantic or affective conflict between text and	266
219	vision (Cai et al., 2019; Pan et al., 2020; Liang et al.,	267
220	2022). Despite strong performance on social media	268
	benchmarks, these methods are limited in visual	269
	symbol grounding and deep metaphor reasoning,	
	both of which are crucial for satirical cartoons and	
	are often cross modal consistent.	
	With the rapid progress of multimodal founda-	
	tion models, recent work leverages LLMs to	
	reformulate MSD. (Tang et al., 2024) uses in-	
	context learning with retrieved similar demonstra-	
	tions, (Jana et al., 2024) applies soft prompt tun-	
	ing to highlight cross modal conflicts, and (Wang	
	et al., 2024b) employs a multiagent pipeline that	
	combines semantic and affective analyses for final	
	prediction. However, these methods have two key	
	limitations. First, they rarely provide verifiable evi-	
	dence and explicit reasoning, making errors hard	
	to diagnose and explanations less reliable. Sec-	
	ond, they inadequately model visual symbols and	
	metaphors.	
	To address this limitation, we propose novel	
	DOUBLE, which consists of Clue2View and SLM	
	Arbiter. It detects satire by measuring the inten-	
	tional divergence between two views, while ex-	
	PLICITLY modeling the underlying reasoning pro-	
	cess, which remains informative even under surface	
	cross modal consistency. DOUBLE further em-	
	ploys a lightweight SLM Arbiter to refine MLLM	
	outputs, mitigating hallucinations while maintain-	
	ing effective and efficient inference.	
	3 MSCE Benchmark	
	In this study, we introduce MSCE, a benchmark of	
	editorial cartoons where satire is conveyed through	
	symbolism, metaphor, and contextual inference	
	rather than explicit image-text contradiction.	
	3.1 Data Collection	
	We construct our benchmark by crawling editorial	
	cartoons from <i>Satire and Humor</i> , a comic supple-	
	ment of <i>People’s Daily</i> . The collected cartoons	
	span a diverse set of topics, including international	
	affairs, domestic social news, and issues related to	
	everyday livelihood. In total, we gather 2,016 car-	
	toon instances. Each instance contains a title and	
	the cartoon image, and a subset further includes	
	editor-provided background descriptions released	
	by <i>Satire and Humor</i> .	
	For samples whose background descriptions are	
	excessively long, we manually rewrite the descrip-	
	tions to improve conciseness and consistency, en-	
	suring that each rewritten background context is	
	limited to approximately 100 Chinese characters	

Text Input	Class	Train	Validation	Test	Total
Title	Sarcastic	801 (73.49%)	44 (70.97%)	91 (71.65%)	1,279 (63.44%)
	Non-Sarcastic	289 (26.51%)	18 (29.03%)	36 (28.35%)	
Title + Background	Sarcastic	419 (67.26%)	25 (65.79%)	55 (72.37%)	737 (36.56%)
	Non-Sarcastic	204 (32.74%)	13 (34.21%)	21 (27.63%)	
Total	Sarcastic	1,220 (71.22%)	69 (69.00%)	146 (71.92%)	2,016
	Non-Sarcastic	493 (28.78%)	31 (31.00%)	57 (28.08%)	

Table 1: Distribution of satiric and non-satiric samples across dataset splits under different text-input settings.

while preserving the essential context required to interpret the satire.

3.2 Dataset Annotation

To ensure high-quality annotations, we conduct manual labeling using the open-source platform Label Studio (Tkachenko et al., 2020–2025). Three master’s students conducted three rounds of independent annotation and verification. All annotators are native Chinese speakers and possess sufficient linguistic proficiency and reading comprehension to reliably identify satirical intent in cartoons.

Prior to annotation, we randomized the distribution of samples to ensure each annotator received a diverse subset of cartoons. During the labeling process, annotators were presented with either (i) the cartoon image and its title, or (ii) the image, title, and accompanying background context. Ground truth labels were determined via majority vote. Detailed annotation instructions can be found in Appendix A. We assess inter-annotator reliability using Fleiss’ Kappa and obtain a score of 0.6302, indicating substantial agreement among annotators (Landis and Koch, 1977).

3.3 Dataset Analysis

Table 1 reports the class distribution across all splits. Overall, (1) the benchmark contains 2,016 instances, with 1,713/100/203 samples in train/validation/test, respectively. (2) Most samples provide only the title as textual input (63.44%). (3) Sarcastic cartoons dominate all splits (69%–72%), indicating a stable class distribution.

Figure 2 reports topic-level statistics of our dataset. Figure 2 (a) shows the relative prevalence of 11 topics, with Economy accounting for the largest share (15.8%), followed by Workplace (14.5%) and Education (14.1%). Technology (12.6%) and Environment (11.5%) also constitute substantial portions, whereas Real Estate is the least represented topic (0.7%). Figure 2 (b)

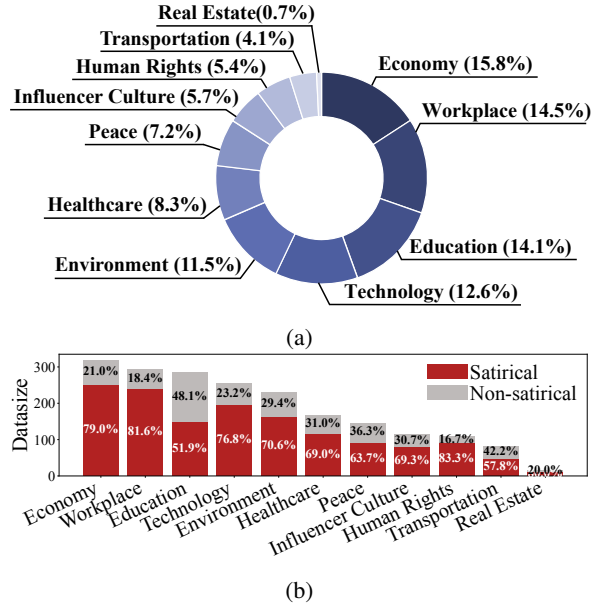


Figure 2: Dataset statistics by topic: (a) topic distribution and (b) satire distribution within each topic.

further breaks down each topic into sarcastic and non-sarcastic subsets. Satire is dominant in most topics, while Education is nearly balanced between the two classes.

4 Methodology: DOUBLE

Problem Formulation. Let $\mathcal{C} = \{\mathcal{T}, \mathcal{V}\}$ denote an editorial cartoon instance, comprising a textual modality \mathcal{T} and a visual modality \mathcal{V} . The goal of MSD is to determine whether \mathcal{C} conveys *satire* by jointly modeling both modalities.

Overview. We first propose Clue2View, a dual-view explanation mechanism for MSD. It explicitly models interpretation through two complementary views. (1) On-Canvas View guides MLLMs to parse the cartoon and text into a Literal Scene Graph (LSG), which captures explicitly depicted entities and their spatial relations, providing a faithful description of the visual content. Based on the LSG, the model then produces a *non-satirical*

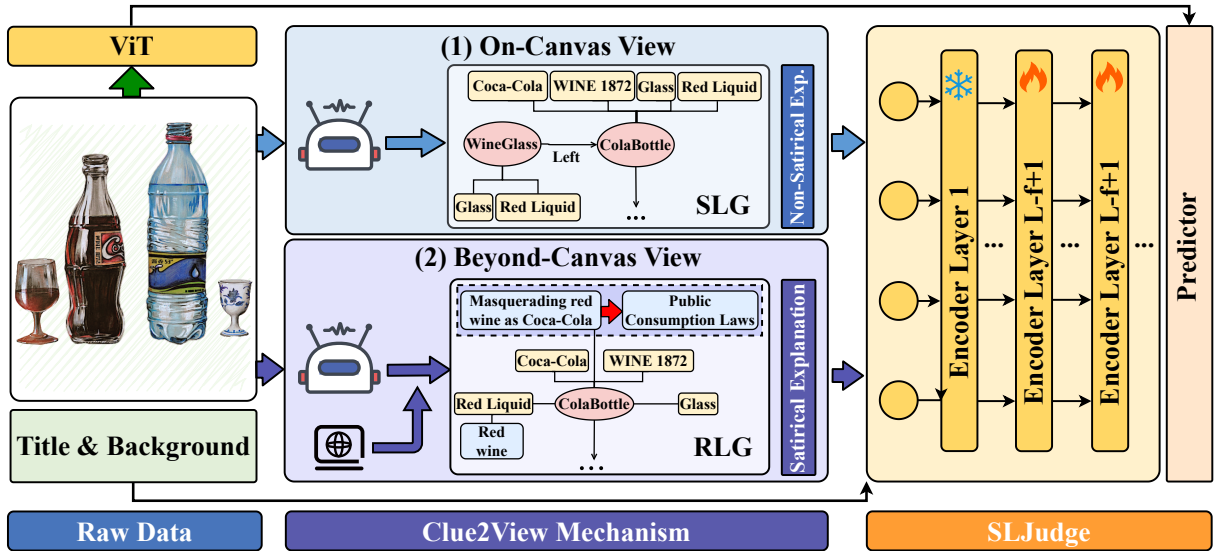


Figure 3: The structure of our proposed DOUBLE framework. (1) Clue2View provides complementary literal and satirical interpretations based on the Literal and Refined Scene Graphs, respectively. (2) The SLM Arbitrator is to distill the reasoning traces from Clue2View to facilitate reliable satire detection.

interpretation that summarizes the literal, factual content. (2) Beyond-Canvas View further instructs MLLMs to enrich the LSG with commonsense and world knowledge, yielding a Refined Scene Graph (RSG) that explicitly represents intent conflicts underlying satire. Conditioned on the RSG, the model generates a *satirical* interpretation that explains the implied meaning. We further introduce SLM Arbitrator to distill and validate the reasoning traces from MLLMs, improving the reliability of the final prediction. Figure 3 provides an overview of our proposed DOUBLE framework.

4.1 Clue2View Mechanism

For the textual modality \mathcal{T} , we extract the cartoon metadata, including the title and the background context (when available), and denote it as \mathcal{T}_m . For the visual modality, we encode the cartoon image and obtain its visual representation, denoted as \mathcal{I}_c . The detailed feature extraction procedure is provided in Appendix D.

4.1.1 On-Canvas View

Cartoonists often conceal their critical stance behind seemingly ordinary scenes, intentionally creating a gap between the literal depiction and the underlying intention. Existing multimodal models often rely primarily on shallow text image alignment cues, rather than high level evidence essential for satire such as intent conflicts and implied stances. Consequently, it is necessary to reorganize raw inputs into structured, reasoning-oriented clues.

To address these problems, Clue2View first refines the multimodal content by extracting salient visual and textual clues to construct a multimodal scene graph, and then rewrites the graph into semantically grounded statements.

We first develop a vision-centric strategy to refine the literal scene content of multimodal comics. The strategy retains only information that is directly observable in the image and the accompanying text. Specifically, it preserves explicit entities, relations, and attributes while discarding external world knowledge, cultural assumptions and information that requires metaphorical interpretation. This restriction yields an objective and intuitive representation that is grounded in the comic itself.

Given a multimodal comic instance \mathcal{C} , we generate a literal scene graph \mathcal{L}_{sg} following the template prompt Prompt_1 in Step 1 of Figure 6 in Appendix D. The generation process is:

$$\mathcal{L}_{sg} = \mathcal{F}(\mathcal{C}, \text{Prompt}_1), \quad (1)$$

where $\mathcal{F}(\cdot)$ denotes the MLLM.

Moreover, since a non-sarcastic reading is essentially an objective statement of the explicit content conveyed by the sample, we further derive a non-sarcastic interpretation from the \mathcal{L}_{sg} . Specifically, guided by the template prompt Prompt_2 in Step 2 of Figure 6 in Appendix D, we obtain the non-sarcastic interpretation: \mathcal{V}_{non} , defined as:

$$\mathcal{V}_{non} = \mathcal{F}(\mathcal{L}_{sg}, \text{Prompt}_2). \quad (2)$$

4.1.2 Beyond-Canvas Sarcastic Interpretation

This step guides the MLLMs to construct a refined scene graph \mathcal{R}_{sg} for a given multimodal satirical comic. We construct \mathcal{R}_{sg} from the \mathcal{L}_{sg} and a fact-oriented description \mathcal{D} . The description serves as a high-level semantic anchor, which specifies the retrieval intent and constrains the integration of relevant world knowledge. In other words, this design reduces semantic drift during knowledge injection and emphasizes the satirical focus. As a result, the \mathcal{R}_{sg} integrates external commonsense knowledge with the information in \mathcal{D} , which improves structured modeling of satirical semantics. For a multimodal comic sample \mathcal{C} , the \mathcal{R}_{sg} is generated using the template Prompt_3 in Step 3 of Figure 6 in Appendix D. We formalize this process as:

$$\mathcal{R}_{sg} = \mathcal{F}([\mathcal{L}_{sg}; \mathcal{D}], \text{Prompt}_3), \quad (3)$$

where $[\cdot; \cdot]$ denotes concatenation.

Based on the semantic relations and metaphorical cues provided by the \mathcal{L}_{sg} and the \mathcal{R}_{sg} , the model can further perform a deeper sarcastic interpretation and inference. The model identifies key semantic conflicts and forms explicit reasoning traces. These traces explain the satirical interpretation, thereby enhancing both interpretability and transparency. The sarcastic reasoning output is guided by the template Prompt_4 in Step 4 of Figure 6 in Appendix D, and is formalized as:

$$\mathcal{V}_{satire} = \mathcal{F}([\mathcal{L}_{sg}; \mathcal{R}_{sg}], \text{Prompt}_4). \quad (4)$$

By leveraging the logically rigorous and coherent Clue2View mechanism, our framework effectively addresses key limitations of prior work (Wang et al., 2024b; Tang et al., 2024; Jana et al., 2024). Existing MSD methods often rely primarily on cross-modal incongruity, overlook hallucinations in MLLMs predictions, and provide limited semantic analysis for metaphor-rich cartoons, which leads to weak performance on non-obvious satire. In contrast, Clue2View explicitly models both literal semantics and deeper intentions, and derives complementary reasoning traces from these representations. This explicit reasoning modeling yields more informative evidence for decision making and translates into substantial gains in prediction accuracy.

4.2 Small Language Model Arbiter

Although MLLMs can produce predictions through our carefully designed four-step pipeline, satire

detection fundamentally relies on holistic and intuitive cognition rather than a sequence of logical inferences (Zhang et al., 2025). This characteristic makes it difficult for MLLMs to reach stable and accurate judgments. Moreover, MLLMs suffer from systematic unfaithfulness (Turpin et al., 2023): their intermediate reasoning and generated rationales may not faithfully reflect the actual decision process, and can even contain hallucination. While task-specific fine-tuning of MLLMs for MSD can partially mitigate these issues, the practical utility of this approach is limited by the massive parameter count of the MLLMs and the corresponding computational costs.

To address these limitations, we propose SLM Arbiter, which distills the outputs of MLLM into a substantially SLM. By fine-tuning this compact model, we obtain a more reliable and deployable solution with far fewer trainable parameters and markedly improved computational efficiency. The fine-tuning procedure can be formulated as:

$$\mathcal{S}(x) = \mathcal{L}_L \circ \mathcal{L}_{L-1} \circ \dots \circ \mathcal{L}_{L-f+1} \circ \mathcal{L}_{L-f}^{\text{fixed}} \circ \dots \circ \mathcal{L}_1^{\text{fixed}}(x), \quad (5)$$

where \mathcal{L}_i denotes the i -th Transformer encoder layer. During training, only the final f layers, from \mathcal{L}_{L-f+1} to \mathcal{L}_L , are updated, while all earlier layers remain frozen. This design preserves the capacity of SLM Arbiter to audit and verify the outputs of the MLLM, while achieving significantly lower training and inference costs.

We concatenate the text \mathcal{T}_m from a satirical cartoon with a non-satirical interpretation \mathcal{V}_{non} and a satirical interpretation \mathcal{V}_{satire} , and feed the resulting sequence into SLM Arbiter to obtain the textual feature representation required for prediction:

$$\mathbf{H}_t = \mathcal{S}([\mathcal{T}_m; \mathcal{V}_{non}; \mathcal{V}_{satire}]), \quad (6)$$

where n denotes the concatenated sequence length and d_t is the hidden dimension. For the visual representation $\mathbf{H}_v \in \mathbb{R}^{n_v \times d_v}$, we extract features from the cartoon image \mathcal{I}_c using a pretrained Vision Transformer (ViT) (Zhang and Yang, 2021): $\mathbf{H}_v = \text{ViT}(\mathcal{I}_c)$, where n_v and d_v denote the numbers of patches and the feature dimension, respectively.

To align the text feature \mathbf{H}_t and the visual feature \mathbf{H}_v , we use MLP-based nonlinear projections $\Psi_t(\cdot)$ and $\Psi_v(\cdot)$ to map them into a shared space. We fuse the projected features with an attention-weighted sum and feed the fused vector to an MLP predictor to obtain $\hat{y} = \text{Predictor}(\text{Fusion}(\Psi_t(\mathbf{H}_t), \Psi_v(\mathbf{H}_v)))$. During training, we optimize the parameters of SLM

Methods	ACC	M-F1	M-P	M-R
Docmsu	72.41	65.85	65.85	65.85
DAIE	81.77	72.94	82.42	70.22
MMSD	67.49	60.96	60.71	61.36
MMSD2.0	76.85	68.17	71.43	66.79
S ³ Agent+GPT-4o-mini	71.92	67.62	66.98	69.25
S ³ Agent+GPT-5.1	82.27	76.12	79.39	74.30
S ³ Agent+Qwen2.5-VL	73.89	65.95	67.08	65.27
COC+GPT-4o-mini	66.90	65.92	68.96	73.31
COC+GPT-5.1	70.44	69.31	71.69	76.78
COC+Qwen2.5-VL	78.82	75.15	74.22	76.72
GOC+GPT-4o-mini	74.88	66.39	68.34	65.42
GOC+GPT-5.1	79.31	75.61	74.71	77.06
GOC+Qwen2.5-VL	76.35	61.15	75.16	60.57
GPT-4o-mini	80.30	74.47	76.22	72.40
GPT-5.1	75.86	72.16	71.29	74.13
Qwen2.5-VL	81.28	77.07	76.80	77.36
DOUBLE+GPT-4o-mini	85.22	81.90	81.58	82.24
DOUBLE+GPT-5.1	84.73	81.57	80.87	<u>82.43</u>
DOUBLE+Qwen2.5-VL	85.22	82.26	81.43	83.31

Table 2: Performance comparison on the MSCE dataset. The best results are highlighted in bold, while the second-best results are indicated with underlines. Higher values of ACC, M-F1, M-P, and M-R indicate better performance.

Arbiter by minimizing the binary cross-entropy loss.

5 Experiments

5.1 Experimental Setup

In this section, we present a simple experimental setup, and the detailed version can be found in the Appendix C.

Baselines. To assess the effectiveness of DOUBLE, we compare it with several competitive baselines in MSD, including SarcasmCue (Yao et al., 2025), S³ Agent (Wang et al., 2024b); Docmsu (Du et al., 2024), DAIE (Wu et al., 2025), MMSD (Cai et al., 2019), MMSD2.0 (Qin et al., 2023).

Metrics. Following prior works (Yao et al., 2025; Hong et al., 2025), we evaluate all models on our proposed MSCE benchmark and employ four metrics : Accuracy (**ACC**), Macro-F1 scores (**M-F1**), Macro Precision (**M-P**) and Macro Recall (**M-R**).

5.2 Main Performance

To assess the effectiveness of DOUBLE, we compare it with 10 competitive baselines, with results reported in Table 2. We draw four key observations.

First, DOUBLE achieves the best overall performance on MSCE across nearly all metrics, improving Accuracy by 12.62% and Macro F1 by 18.06% over the average baseline. The advantage remains

consistent when switching backbones, indicating robust gains rather than backbone dependent tuning.

Second, methods relying on shallow cross modal cues are insufficient for handling the deep metaphorical expressions that frequently appear in MSCE. Such approaches typically emphasize surface level cross modal inconsistency. However, the sarcastic intent is often implicit in MSD, the metaphor is deeper, and cross modal conflict may not be explicitly manifested. As a result, these models struggle to capture the core semantics, which fundamentally limits their overall performance.

Third, longer and more complex reasoning chains do not necessarily help. Multiagent reasoning and chain based approaches such as S³ Agent, COC, and GOC show unstable gains and even notable degradation on some backbones, suggesting that extended reasoning is vulnerable to hallucination and noise accumulation that corrupts crucial evidence.

Finally, while MLLMs based methods are strong in zero-shot multimodal settings, they are less reliable on MSD without task specific adaptation, and full fine tuning is often prohibitively expensive. By refining raw the outputs of MLLMs through SLM Arbiter, DOUBLE provides more stable predictions and achieves better MSD performance than conventional MLLMs based approaches.

5.3 Ablation Study

Effect of Clue2View Mechanism. To quantify the contribution of each component in Clue2View mechanism, we conduct an ablation study that isolates individual steps and measures their impact on overall performance. Specifically, we evaluate four variants: (1) *w/o LSG*, which removes the LSG extraction module; (2) *w/o RSG*, which removes the RSG extraction module; (3) *w/o SG*, which removes both LSG and RSG and disables the dual view interpretation mechanism accordingly; and (4) *w/o Dual View*, which removes dual view interpretation and models the task using only scene graph description.

The results indicate that each step of Clue2View is essential for satirical cartoon detection. Beyond improving accuracy, these components also provide useful interpretability by offering reasoning evidence that supports the final prediction. Notably, removing the RSG mechanism yields the largest performance drop. This suggests that satirical cartoons often convey intent through implicit mean-

Module	MLLM Backbone	Variant	ACC	M-F1
Clue2View	GPT-4o-mini	w/o LSG	84.73	80.33
		w/o RSG	82.76	78.54
		w/o SG	84.24	79.05
		w/o Dual-View	84.73	81.57
Clue2View	GPT-5.1	w/o LSG	84.24	79.57
		w/o RSG	83.74	79.54
		w/o SG	84.24	79.05
		w/o Dual-View	84.24	80.89
Clue2View	Qwen2.5-VL	w/o LSG	83.74	79.06
		w/o RSG	83.25	79.04
		w/o SG	84.24	79.05
		w/o Dual-View	85.22	81.71
SLM Arbiter	GPT-4o-mini	w/o Fine-tune	79.31	73.81
	GPT-5.1	w/o Fine-tune	79.31	73.51
	Qwen2.5-VL	w/o Fine-tune	78.82	73.04
DOUBLE	GPT-4o-mini	ALL	85.22	81.90
	GPT-5.1	ALL	84.73	81.57
	Qwen2.5-VL	ALL	85.22	82.26

Table 3: Ablation study on key components of DOUBLE.

ings and deeper metaphors, where relying solely on surface level clues is insufficient to consistently capture the key semantics. In contrast, the deep semantic reasoning cues provided by RSG are more critical for robust discrimination.

Effect of SLM Arbiter. To assess the effectiveness of SLM Arbiter, we introduce an ablation variant *w/o Fine tune*, where the SLM is not fine tuned. The results show that fine tuning SLM Arbiter is necessary to adapt to the distilled outputs produced by the MLLM. Freezing all layers leads to a substantial performance degradation, confirming the importance of parameter updating for effective alignment with the distillation signals.

5.4 Case Study

To further examine the interpretability of our proposed DOUBLE, we randomly sample a cartoon from the MSCE dataset and analyze how DOUBLE infers its satirical label, as shown in Figure 4. The cartoon depicts several drink containers and each associated with different beverages. DOUBLE first constructs a LSG and produces an intuitive non-satirical reading, namely that the image portrays a gathering that prepares both alcoholic and non alcoholic drinks to accommodate different preferences. It then incorporates relevant commonsense and contextual knowledge, such as the Central Eight Point Regulations and related laws, to refine the scene representation. Based on the resulting RSG, DOUBLE performs further reasoning and derives the underlying satirical interpretation, which uses the surface narrative of 'private customization' to criti-

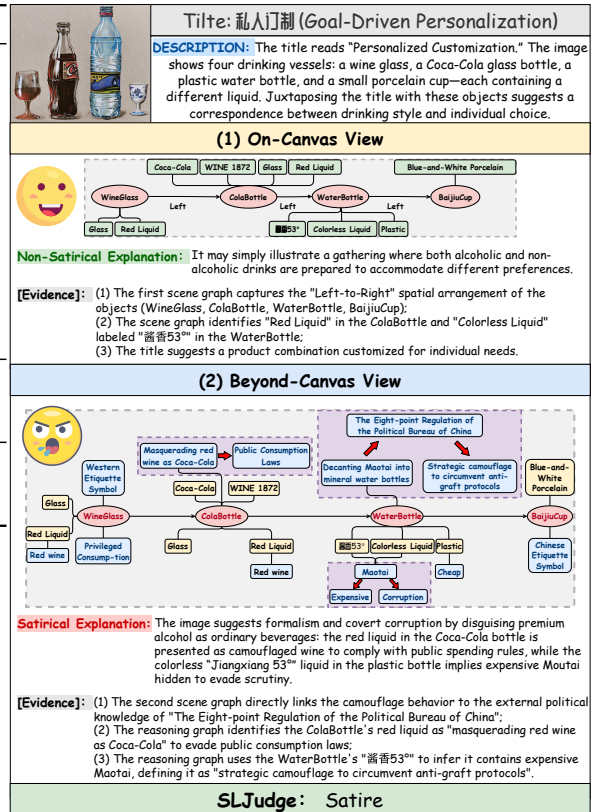


Figure 4: Case study of correctly detected comic satire on the MSCE dataset.

cize deliberate rule breaking and concealed corruption.

This case illustrates that DOUBLE not only detects satirical content, but also provides a coherent and well supported justification grounded in commonsense knowledge and logical inference.

6 Conclusion

This work introduces MSCE, the first official media sourced benchmark for MSD. MSCE emphasizes hard cases where image text pairs appear coherent, but satire arises from metaphor, symbolism, and culture specific context, demanding nontrivial cross modal reasoning. We also propose DOUBLE, an explainable framework for MSD. Through Clue2View, DOUBLE progressively contrasts a literal reading with the implied intent. A Small Language Model Arbiter further distills the dual view traces from the MLLM to support efficient and reliable decisions, reducing cost and filtering hallucinated evidence. Experiments and case studies show that DOUBLE achieves strong performance with clear rationales, providing a competitive baseline for MSD.

7 Limitations

Cross-cultural and time-sensitive generalization.

Satire interpretation is strongly conditioned on cultural context and current affairs. Even with Beyond-Canvas knowledge enrichment, cartoons from different regions and languages can differ substantially in symbolism, sociopolitical framing, metaphor conventions, and textual expression, which may reduce robustness under domain shift. Moreover, DOUBLE relies on the coverage and freshness of the underlying MLLMs knowledge; for newly emerging events or memes, the model may lag or misinterpret. Biases in long-tail or niche knowledge can further affect the construction of the refined scene graph and the induced intent conflicts.

Subjectivity and annotator agreement.

Satire is inherently subjective, and judgments vary with annotators' background knowledge, cultural assumptions, and personal stance. Despite using standardized guidelines and consistency checks, disagreements and borderline cases are difficult to eliminate. This residual label noise can blur decision boundaries and make evaluation sensitive to small changes in annotation criteria.

8 Ethical considerations

We have obtained the necessary authorization to collect and use content from People's Daily Online, as the website's publicly available terms grant rights to users who comply with its online rules. We carefully reviewed these terms and adhered strictly to all stated requirements throughout data acquisition and annotation. In the revised manuscript, we provide a detailed account of the data sources, the collection pipeline, and the annotation procedure to ensure transparency and reproducibility. We further commit to continuous oversight of the dataset and will promptly remove or revise any data that could be considered ethically inappropriate.

This work advances multimodal sarcasm detection to enable responsible research use and to deepen our understanding of sarcastic expression in real-world communication, thereby supporting content analysis and related downstream applications. We acknowledge that techniques developed for sarcasm detection may be misused to facilitate the creation of misleading content that is more difficult to identify. We explicitly oppose any such misuse. To mitigate this risk, we will release the

dataset under an online usage agreement that requires all users to comply with applicable website terms and relevant legal and ethical standards, and that prohibits use intended to produce deceptive or harmful content.

All annotators of our dataset are coauthors of this manuscript rather than crowdworkers. They participated as invited volunteers (without using any third-party crowdsourcing platform), and the annotation process was conducted under our direct supervision.

References

- Michele Bedard and Chianna Schoenthaler. 2018. Satire or fake news: Social media consumers' socio-demographics decide. In *Companion proceedings of the the web conference 2018*, pages 613–619.
- Yitao Cai, Huiyu Cai, and Xiaojun Wan. 2019. Multimodal sarcasm detection in twitter with hierarchical fusion model. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 2506–2515.
- Sohan De Sarkar, Fan Yang, and Arjun Mukherjee. 2018. Attending sentences to detect satirical fake news. In *Proceedings of the 27th international conference on computational linguistics*, pages 3371–3380.
- Poorav Desai, Tanmoy Chakraborty, and Md Shad Akhtar. 2022. Nice perfume. how long did you marinate in it? multimodal sarcasm explanation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10563–10571.
- Hang Du, Guoshun Nan, Sicheng Zhang, Binzhu Xie, Junrui Xu, Hehe Fan, Qimei Cui, Xiaofeng Tao, and Xudong Jiang. 2024. Docmsu: A comprehensive benchmark for document-level multimodal sarcasm understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17933–17941.
- Shafkat Farabi, Tharindu Ranasinghe, Diptesh Kanojia, Yu Kong, and Marcos Zampieri. 2024. A survey of multimodal sarcasm detection. In *IJCAI*.
- Rongpei Hong, Jian Lang, Jin Xu, Zhangtao Cheng, Ting Zhong, and Fan Zhou. 2025. Following clues, approaching the truth: Explainable micro-video rumor detection via chain-of-thought reasoning. In *Proceedings of the ACM on Web Conference 2025*, pages 4684–4698.
- Soumyadeep Jana, Animesh Dey, and Ranbir Singh Sanasam. 2024. Continuous attentive multimodal prompt tuning for few-shot multimodal sarcasm detection. In *Proceedings of the 28th Conference on Computational Natural Language Learning*, pages 314–326.

827 An Yang, Junshu Pan, Junyang Lin, Rui Men, Yichang 876
828 Zhang, Jingren Zhou, and Chang Zhou. 2022. Chi- 877
829 nese clip: Contrastive vision-language pretraining in 878
830 chinese. *arXiv preprint arXiv:2211.01335*.
831 Ben Yao, Yazhou Zhang, Qiuchi Li, and Jing Qin. 2025. 879
832 Is sarcasm detection a step-by-step reasoning pro- 880
833 cess in large language models? In *Proceedings of* 881
834 *the AAAI Conference on Artificial Intelligence*, vol- 882
835 ume 39, pages 25651–25659. 883
836 Tan Yue, Xuzhao Shi, Rui Mao, Zonghai Hu, and Erik 884
837 Cambria. 2024. Sarcnet: a multilingual multimodal 885
838 sarcasm detection dataset. In *Proceedings of the* 886
839 *2024 Joint International Conference on Computa-* 887
840 *tional Linguistics, Language Resources and Evalua-* 888
841 *tion (LREC-COLING 2024)*, pages 14325–14335. 889
842 Qinglong Zhang and Yu-Bin Yang. 2021. Rest: An effi- 890
843 cient transformer for visual recognition. volume 34, 891
844 pages 15475–15485. 892
845 Yazhou Zhang, Chunwang Zou, Zheng Lian, Prayag 893
846 Tiwari, and Jing Qin. 2025. Sarcasmbench: Towards 894
847 evaluating large language models on sarcasm under- 895
848 standing. *IEEE Transactions on Affective Comput-* 896
849 *ing*. 897

850 Appendix 898

851 A Human Annotation Instructions 899

852 This appendix describes the human annotation pro- 900
853 tocol for satirical cartoons, including the opera- 901
854 tional definition of satire, the labeling rationale, 902
855 and the practical workflow implemented with La- 903
856 bel Studio. 904

857 A.1 Operational Definition of Satire 905

858 In this work, *satire* is defined as the use of rhetori- 906
859 cal devices such as metaphor and exaggeration to 907
860 expose, criticize, or ridicule a person, group, or 908
861 social phenomenon. We adopt the definition from 909
862 *Xinhua Dictionary* as the normative reference. For 910
863 annotation, satire is treated as a semantic structure 911
864 in which the surface meaning is intentionally mis- 912
865 aligned with the intended meaning, commonly re- 913
866 alized through irony, contrast, or visually grounded 914
867 critique. 915

868 A.2 Labeling Scheme and Aggregation Rule 916

869 To reduce ambiguity during annotation, we employ 917
870 three labels that explicitly distinguish non satiri- 918
871 cal content from satire with different degrees of 919
872 directedness or hostility. This three class design 920
873 is introduced primarily to make the annotator’s de- 921
874 cision process easier and more consistent. In the 922
875 final dataset used for modeling and evaluation, we 923

876 merge label=1 and label=2 into a single positive 877
878 category, namely *satire*, while label=0 remains 879
880 the negative category. 881

882 When an instance is difficult to judge, annota- 883
884 tors may consult an AI assistant by providing the 884
885 definitions and diagnostic cues in this section as 886
887 contextual information. The assistant’s explanation 888
889 is used only as a reference. The final label is always 890
891 determined by the human annotator. 892

893 A.3 Label Definitions and Diagnostic Cues 898

894 **Label 0: Non Satirical.** An instance is labeled as 895
896 non satirical when neither the image nor the title 897
898 exhibits a satirical structure. In particular, there 899
900 is no discrepancy between literal expression and 901
902 intended meaning, no conflict between the title 903
904 and the depicted content, no visual metaphor that 904
905 implies a directed critique, and no irony. Typical 905
906 instances show high consistency between the title 906
907 and the image, and they do not present an identifi- 907
908 able target of criticism. Exaggeration may appear 908
909 for artistic or humorous effect, and emotional tones 909
910 such as resignation, light humor, or lyricism are 910
911 allowed as long as they do not function as satire. 911

912 **Label 1: Non Aggressive and Non Malicious** 912
913 **Satire.** An instance is labeled as non aggressive 913
914 satire when a satirical structure is present but the 914
915 instance does not attack a person or an entity. The 915
916 intent is primarily self directed or situation oriented, 916
917 such as self deprecating irony or a non confronta- 917
918 tional portrayal of everyday predicaments. The 918
919 instance typically lacks a specific target of attack, 919
920 and its affect often reflects resignation, absurdity, 920
921 or mild complaint rather than hostility. 921

922 **Label 2: Aggressive and Malicious Satire.** An 922
923 instance is labeled as aggressive satire when a satiri- 923
924 cal structure is present and the instance conveys a 924
925 directed intent to criticize, attack, or expose. This 925
926 class is characterized by contradiction, tension, or 926
927 metaphorical mismatch between the title and the 927
928 image, and by a clearly identifiable target of cri- 928
929 tique. Common targets include governments, polit- 929
930 ical systems, and bureaucratic structures; capital, 930
931 corporations, and workplace exploitation; unequal 931
932 allocation of social resources, class stratification, 932
933 wealth inequality, and unequal opportunity; pub- 933
934 lic services such as education, healthcare, elder 934
935 care, and housing; technology, media, and digi- 935
936 tal life including smartphones, the internet, online 936
937 platforms, and advertising; environmental and eco- 937
938 logical issues such as pollution, performative envi- 938
939 ronmentalism, ecological destruction, and resource 939
940

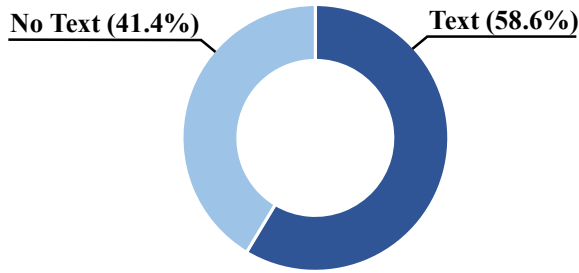


Figure 5: Percentage of images containing embedded text.

waste; war, violence, security, terrorism, public safety, and arms races; and collective social behavior and norms such as bandwagon effects, indifference, moral numbness, and conformity. Targets outside these categories are also allowed when the critique is explicit.

A.4 Annotation Workflow in Label Studio

The annotation environment is reproduced locally using Label Studio. Annotators install Label Studio in a VS Code terminal with Python version at least 3.8 by running `pip install label-studio`. They then start the service by running `label-studio` and keep the terminal session active during annotation. The interface opens at `http://localhost:8080`. Annotators create a local account on first use and log in, noting that all data are stored locally.

A new project is created with the name `sarcasm_detection`. The dataset file `clean_data.json` is uploaded through the `Data Import` entry. The labeling interface is configured by pasting the contents of `templates.txt` into the template editor and saving the configuration. Annotation is performed by entering the project and selecting `Label All Tasks`, assigning exactly one label to each instance according to the definitions and cues above. After all instances are labeled, the results are exported using the built in export function in the Label Studio interface.

B Dataset Analysis

Figure 5 summarizes the prevalence of embedded text in images: 58.6% of images contain on-image text, while 41.4% do not.

C Experiment Setup

C.1 Description of Baselines

In our experiments we compare against both multimodal sarcasm detection methods that are built on

top of multimodal large language models (MLLMs) and those that are *not* based on MLLMs. The distinction is important: methods based on MLLMs exploit large pretrained vision–language reasoning capabilities to perform inference with minimal task-specific training, whereas the non-MLLM baselines typically rely on task-specific feature fusion and dedicated classifiers.

The baselines we evaluate include the following five representative approaches.

- **SarcasmCue** (Yao et al., 2025) investigates whether large language models can detect sarcasm via intermediate reasoning steps. It introduces a prompting framework with multiple strategies that elicit reasoning about linguistic and contextual cues and evaluates LLMs on standard sarcasm benchmarks. This study highlights that non-sequential prompting often outperforms sequential prompting even for LLMs, suggesting that sarcasm comprehension may not strictly follow a linear reasoning process.
- **S³ Agent** (Wang et al., 2024b) proposes a multi-view agent framework that unlocks the power of vision large language models for zero-shot multimodal sarcasm detection. By decomposing sarcasm signals into superficial expression, semantic information, and sentiment views, this method leverages the reasoning ability of pretrained vision LLMs to achieve strong zero-shot performance on MMSD2.0.
- **DocMSU** (Du et al., 2024) introduces a comprehensive benchmark for document-level multimodal sarcasm understanding in news articles. In addition to the dataset, the authors propose a baseline approach that aligns pixel-level image features with word-level textual representations to capture sarcasm cues across long texts and heterogeneous content.
- **Dual-Level Adaptive Incongruity-Enhanced model (DAIE)** (Wu et al., 2025) is a deep architecture emphasizing multimodal incongruity. It encodes text and visual modalities and then fuses information at token and graph levels to explicitly highlight contradiction patterns, which are central to multimodal sarcasm detection.

1012 • **MMSD2.0** (Qin et al., 2023) defines a cor-
 1013 rected multimodal sarcasm detection system
 1014 and benchmark that addresses issues in the
 1015 original MMSD dataset by removing spurious
 1016 cues and reannotating unreasonable samples.
 1017 A multi-view CLIP framework is introduced
 1018 to fuse text, image, and text-image interac-
 1019 tion representations, substantially improving
 1020 reliability and baseline performance on this
 1021 benchmark.

1022 • **Twitter Hierarchical Fusion Model**
 1023 (MMSD) (Cai et al., 2019) represents an
 1024 early multimodal sarcasm detection approach
 1025 on social media, using hierarchical fusion
 1026 of text and image features with attention
 1027 mechanisms. This model integrates modality
 1028 features at multiple levels to highlight
 1029 cross-modal inconsistencies for sarcasm
 1030 detection, though without relying on MLLMs.

1031 Among these baselines, the approaches that are
 1032 based on multimodal large language models are S³
 1033 Agent and SarcasmCue, because both take advan-
 1034 tage of large pretrained vision or language models
 1035 to make predictions with little or no task-specific
 1036 training. SarcasmCue investigates how internal
 1037 reasoning processes in large language models re-
 1038 late to sarcasm understanding, and S³ Agent uses
 1039 vision LLMs to address the multimodal sarcasm
 1040 detection problem in a zero-shot scenario. The
 1041 other baselines such as the alignment baseline
 1042 in DocMSU, the dual-level adaptive incongruity
 1043 model, the MMSD2.0 multi-view CLIP framework,
 1044 and the hierarchical fusion model for Twitter focus
 1045 on customized multimodal representations and fea-
 1046 ture fusion rather than on large model reasoning.
 1047 Taken together, these methods provide a compre-
 1048 hensive set of comparative systems for evaluating
 1049 multimodal sarcasm detection performance.

1050 D Implementation Details

1051 In this section, we present detailed implementation
 1052 specifications for our proposed DOUBLE as well
 1053 as a comprehensive overview of the experimental
 1054 setup.

1055 D.1 Implementation Environment

1056 All experiments are conducted on a system compris-
 1057 ing an Intel(R) Core(TM) i9-14900KF processor,
 1058 equipped with one NVIDIA GeForce RTX 4090 D
 1059 GPU with 24 GB of VRAM, and accompanied by
 1060 64 GB of DRAM.

Table 4: Example of prompt for various perspectives applied in S³ Agent.

Superficial Expression Prompt

Given the following image and text, judge whether there is sarcasm based on the superficial expression. This requires detecting underlying critique through discrepancies between image and text, considering both modalities together.
 text: <text>

Semantic Information Prompt

Given the following image and text, judge whether there is sarcasm based on the semantic information. This prompt focuses on extreme portrayals and metaphorical meanings inferred from the combined image and text.
 text: <text>

Sentiment Expression Prompt

Given the following image and text, judge whether there is sarcasm based on sentiment expression. This requires detecting critical emotions targeted at specific subjects or behaviors by jointly examining the image and text.
 text: <text>

D.2 MLLM Implementation in Clue2View

To ensure the robustness and generalizability of our findings, we conduct all experimental evaluations using a diverse set of three state-of-the-art MLLMs: GPT-4o-mini, GPT-5.1, and Qwen2.5-VL-72B-Instruct. By employing this range of models, we verify the effectiveness of our approach across different architectural designs and resource constraints."

D.3 SLM Implementation in SLM Arbiter

Our SLM is based on a masked self-attention Transformer architecture, i.e., BERT, pretrained through language-visual contrastive learning (Radford et al., 2021). For the visual feature encoding, we leverage the pre-trained Vision Transformer (ViT), keeping its parameters frozen. Specifically, for our benchmarks, we employ OFA-Sys/chinese-clip-vit-large-patch14 (Yang et al., 2022).

D.4 Data Preprocessing

Our SLM is based on a masked self-attention Transformer architecture, i.e., BERT, pretrained through language-visual contrastive learning (Rad-

ford et al., 2021). For the visual feature encoding, we leverage the pre-trained Vision Transformer (ViT), keeping its parameters frozen. Specifically, for our benchmarks, we employ OFA-Sys/chinese-clip-vit-large-patch14 (Yang et al., 2022).

Given an image–text pair X , we first extract its multimodal information from both modalities. For the visual content, we represent the image as I and adopt the vision encoder of OFA-Sys/chinese-clip-vit-large-patch14 (Yang et al., 2022) as the feature extractor. Specifically, we feed I into the CLIP vision transformer and take the output embedding of the [CLS] token as the global image representation, denoted as $\mathbf{h}_v \in \mathbb{R}^{d_v}$, where d_v is the dimensionality of the visual embedding space defined by the model.

For the textual content, we use the image associated title and external background knowledge as the textual input. Specifically, we denote the title as $T_{\text{title}} \in \mathbb{R}^{n_t}$ and the background knowledge as $T_{\text{bg}} \in \mathbb{R}^{n_b}$, where n_t and n_b are the numbers of words (or tokens) in the title and the background text, respectively. We then concatenate them to form the final textual sequence $\tilde{T} = [T_{\text{title}}; T_{\text{bg}}] \in \mathbb{R}^{\tilde{n}}$, where $\tilde{n} = n_t + n_b$. This design enables the model to interpret the image in context by jointly considering the visual content, the title, and the relevant background information.

D.5 Training Configuration

Experimental Setup. We implement our model using the PyTorch framework. For text encoding, we utilize the pre-trained Chinese-CLIP-ViT-Large-Patch14 model to extract textual features. For visual features, we align the input dimensions through a multi-layer perceptron (MLP) projection. We employ the AdamW optimizer for model parameter optimization, configured with a learning rate of 7.7×10^{-5} and a weight decay of 5×10^{-4} . To mitigate the class imbalance issue, we adopt the Focal Loss with γ set to 1.8. The maximum number of training epochs is set to 15, with an early stopping mechanism triggered if the validation accuracy does not improve for 5 consecutive epochs. The batch size is set to 32, and the random seed is fixed at 2026 to ensure reproducibility. For statistical reliability, we report the results based on the best-performing model selected on the validation set. For baseline models, we strictly adhere to the settings specified in their original papers.

Prompt₁: Literal Scene Graph

Role: You are designated as a specialized agent for Multimodal Knowledge Graph construction, focusing on surface-level visual and textual grounding. The extraction process is divided into two sequential phases: (A) Initial Multimodal Parsing and (B) Surface-KG Refinement.

Phase A: Fine-Grained Multimodal Parsing Objective: Perform a comprehensive extraction of entities, attributes, and relationships derived solely from the visual content and the provided title metadata. Methodology: Employ multi-scale visual inspection (zooming or cropping) to resolve fine-grained details, particularly for textual elements, symbols, or small objects. Inference Constraint: All extractions must be strictly grounded in observable evidence. Probabilistic inference is permissible only when strongly supported by visual cues; hallucination of unobserved entities, relations, or events is strictly prohibited.

Phase B: Surface-KG Refinement Objective: Critique the initial extraction to produce a finalized, revised Surface Knowledge Graph. This graph must contain only factual assertions directly verifiable in the image, strictly filtering out over-interpreted semantic abstractions.

Operational Protocols and Constraints:

1. Entity Extraction Protocols: Taxonomy: All entity nodes must be labeled in English. Granularity: Decompose composite objects into sub-entities when distinct interactions (e.g., handling, using) are observed. Textual Entities: Extract all visible text strings as distinct entities. Preserve non-English content exactly as it appears within attribute values, but maintain English node labels.
2. Relation Extraction Protocols: Predicate Ontology: Relations must be expressed as specific, English action verbs or spatial prepositions. Grounded Triples: Generate only factual triplets (Entity-Relation-Entity) substantiated by observable spatial arrangements, actions, interactions, or communicative gestures. Speaker-Utterance Alignment: For dialogue or text associated with a specific character, enforce the structure: [Person_Entity] -> [say] -> [Text_Entity].
3. Attribute Extraction Protocols: Visual Attributes: Annotate entities with visually grounded descriptors including color, material, pose, facial expression, relative size, and spatial coordinates. Textual Attributes: For on-screen text, include attributes such as raw_text_content, layout_structure, reading_direction, and text_type.
4. OCR and Layout Analysis (Chinese/Multilingual): Segmentation: Extract visible text fragments as individual entities. Reconstruction: Merge text fragments only when the layout strongly implies continuity; otherwise, treat them as discrete. If the reading order is ambiguous, annotate the direction as 'unknown'.
5. Evidence Modality Annotation: Requirement: Every relation and attribute in the final output must include an 'evidence_modality' tag to specify the information source. Valid Tags: 'image' (visual only), 'text' (textual only), or 'image+text' (cross-modal inference).

Input: Image and Title. **Output:** A structured, purified Surface Knowledge Graph adhering to the above constraints.

(a) Step 1

Prompt₂: Non-Sarcastic Interpretation

You are an expert in abductive reasoning. Your task is to analyze the provided {surface_kg} to generate a streamlined rationale from a literal perspective.

Strictly adhere to the provided Knowledge Graph nodes and edges without hallucinating external information.

Task:

Using ONLY the {surface_kg}, explain why the sample appears reasonable or non-sarcastic in its literal sense. Focus on the visible facts and standard interpretations of the objects and actions described.

Please answer in Chinese.

Return ONLY valid JSON with the following structure:

```
{
  "non_sarcastic_explanation": "string"
}
```

(b) Step 2

Prompt₃: Refined Scene Graph

Role: You are designated as a specialized agent for the construction of Deep Semantic Knowledge Graphs (Deep-KG), specifically tailored for the interpretation of satirical cartoons. Your objective is to model the latent semantic layer—symbolism, intent, and causal logic—essential for computational sarcasm detection.

Input Specifications:

Meta-Information: Title text and optional background context.

Visual Narrative: A concise LLM-generated description of the scene.

Surface Knowledge Graph: A structured representation of explicitly visible facts.

Methodological Framework:

Semantic Abstraction and Inference: Construct a Deep-KG that captures implied meanings, metaphorical mappings, social roles, causal dependencies, communicative intents, normative expectations, and abstract concepts. The graph must bridge the gap between visual signals and high-level interpretation.

Redundancy Constraint: Strictly exclude purely visible facts already present in the Surface KG, unless they serve as necessary semantic anchors for connecting abstract nodes (e.g., grounding a symbol to its physical representation).

Node Ontology and Naming: You are authorized to introduce abstract nodes where supported by the title, background context, or scene composition. All node identifiers must adhere to snake_case formatting. Existing node IDs should be preserved where applicable; new nodes should be generated only to represent distinct semantic concepts.

Output Specifications: Return the result strictly as a valid JSON object adhering to the schema defined below. Do not include markdown formatting, code block delimiters, or explanatory text.

(c) Step 3

Prompt₄: Sarcastic Interpretation

You are an expert in abductive reasoning for sarcasm detection. Your task is to analyze the provided {surface_kg} and {deep_kg} to generate a streamlined rationale from a sarcastic perspective. Strictly adhere to the provided Knowledge Graph nodes and edges without hallucinating external information.

Task:

Combining the {surface_kg} and {deep_kg}, explain why the sample is reasoned as sarcastic. Identify semantic conflicts, common sense violations, or contradictions between the surface layer (what is seen) and the deep layer (symbolism, intent, or cultural context).

Please answer in Chinese.

Return ONLY valid JSON with the following structure: { "sarcastic_explanation": "string" }

(d) Step 4

Figure 6: Prompting overview.