
CAF-I: A Collaborative Multi-Agent Framework for Enhanced Irony Detection with Large Language Models

Ziqi Liu ^{* 1} Ziyang Zhou ^{* 1} Mingxuan Hu ¹

Abstract

Large language model (LLM) have become mainstream methods in the field of sarcasm detection. However, existing LLM methods face challenges in irony detection, including: **1. single-perspective limitations**, **2. insufficient comprehensive understanding**, and **3. lack of interpretability**. This paper introduces the Collaborative Agent Framework for Irony (CAF-I), an LLM-driven multi-agent system designed to overcome these issues. CAF-I employs specialized agents for Context, Semantics, and Rhetoric, which perform multidimensional analysis and engage in interactive collaborative optimization. A Decision Agent then consolidates these perspectives, with a Refinement Evaluator Agent providing conditional feedback for optimization. Experiments on benchmark datasets establish CAF-I’s state-of-the-art zero-shot performance. Achieving SOTA on the vast majority of metrics, CAF-I reaches an average **Macro-F1 of 76.31%**, a **4.98%** absolute improvement over the strongest prior baseline. This success is attained by its effective simulation of human-like multi-perspective analysis, enhancing detection accuracy and interpretability.

1. Introduction

The rapid development of social media has fostered increasingly diverse and semantically complex forms of textual expression (Kader et al., 2022) highlights the challenge of irony detection. Irony, a rhetorical device contrasting stated and intended meaning (Booth, 1974), is crucial for NLP sys-

tems as its presence alters textual interpretation, impacting applications like sentiment analysis and content moderation (Reyes et al., 2013). However, detection is complicated by context dependency, obscure authorial intent, and varied rhetorical strategies (Wallace et al., 2015).

Initial efforts used traditional machine learning with hand-crafted features (Reyes et al., 2013), often failing to capture irony’s subtleties. Deep learning methods improved semantic representation but struggled with the implicit and contextual nature of ironic utterances (Ghosh & Veale, 2016; Baziotis et al., 2018). Consequently, the emergence of LLM has presented novel avenues, exhibiting substantial potential in irony detection, typically via fine-tuning or advanced prompt engineering techniques (Yao et al., 2024).

Despite this significant promise and their advanced capabilities, LLM encounters a distinct set of challenges when specifically applied to the nuanced task of irony detection, as illustrated in Figure 1. Key among these challenges are:

- 1. Single-Perspective Limitation:** As single-model predictors, LLM have limited capacity for multidimensional collaborative reasoning, struggling to collaboratively synthesize diverse analytical insights required for deconstructing complex irony.
- 2. Insufficient Comprehensive Understanding:** LLM struggle to holistically integrate diverse informational cues, such as context, semantics, and rhetoric, to achieve human-like deep comprehension of ironic intent.
- 3. Lack of Interpretability:** Opaque reasoning process in irony judgments hinders the understanding, trust, and debugging of decision-making mechanisms.

Motivated by Chain-of-Thought (CoT) prompting (Wei et al., 2022) and multi-agent cooperation principles (Liu et al., 2023), this paper introduces CAF-I, a novel LLM-driven multi-agent framework to address existing limitations in irony detection by simulating human-like multidimensional analytical processes. CAF-I initially performs multi-perspective parallel analysis from three core angles: contextual consistency, deep semantic logic, and rhetorical function. Subsequently, an interactive collaborative

^{*}Equal contribution ¹School of Advanced Technology, Xi’an Jiaotong-Liverpool University, Location, China. Correspondence to: Ziqi Liu <Ziqi.Liu22@student.xjtlu.edu.cn>, Ziyang Zhou <Ziyang.Zhou22@student.xjtlu.edu.cn>, Mingxuan Hu <Mingxuan.Hu22@student.xjtlu.edu.cn>.

Accepted at the ICML 2025 Workshop on Collaborative and Federated Agentic Workflows (CFAgentic@ICML’25), Vancouver, Canada. July 19, 2025. Copyright 2025 by the author(s).

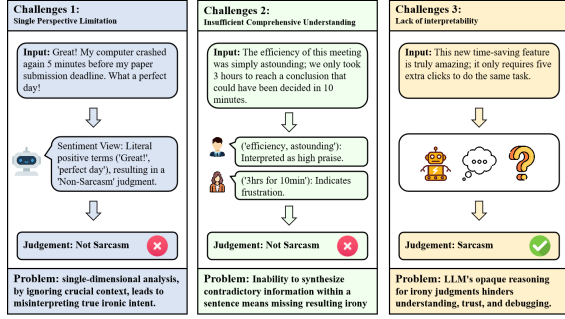


Figure 1. Examples illustrating three primary challenges for LLM in sarcasm detection: single perspective limitation, insufficient comprehensive understanding, and lack of interpretability.

optimization mechanism facilitates information sharing and debate among these perspectives. Finally, an integrated decision mechanism allows a central decision-making agent to consolidate these viewpoints and apply conditional optimization triggered by evaluative feedback, forming the paper’s primary contributions:

- Introduces CAF-I, a novel LLM-driven multi-agent framework significantly improving irony detection through specialized agent collaboration and phased optimization.
- Achieves state-of-the-art zero-shot performance across multiple benchmark datasets.
- Its unique architecture offers clear decision pathways and error correction capabilities, enhancing interpretability.
- Extensive experiments validate its component necessity, overall robustness, and innovation, establishing CAF-I as a pioneering application of specialized LLM agents in structured collaboration for irony detection.

2. RELATED WORK

2.1. Irony Detection

Irony detection in NLP remains challenging due to its inherent complexity, context dependence, and contradictory nature. Early machine learning relied on hand-crafted features, including lexical cues (Davidov et al., 2010) and sentiment lexicons (Reyes et al., 2013), but inadequately captured irony’s subtleties. Deep learning techniques offered significant improvements: Word embeddings such as Word2Vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014) enabled richer semantic representations. Subsequently, Convolutional Neural Networks (Poria et al., 2016), recurrent neural networks, and particularly LSTM architectures (Zhang et al., 2016) learned hierarchical and sequential features, while Graph Neural Networks (Liang et al., 2022)

captured structural dependencies to improve performance. Nonetheless, these deep learning approaches still faced difficulties with irony’s implicit meanings and complex rhetoric, motivating the exploration of more advanced models.

Recent advances in large language models have introduced prompt-based and zero-shot learning methods, allowing models like GPT-4o to interpret irony with minimal reliance on labeled data (Yao et al., 2024). By leveraging task-specific prompts, these models have demonstrated enhanced context-aware irony classification. However, existing LLM-based approaches still operate as single-model predictors, lacking the capacity for multi-perspective reasoning.

2.2. CoT Prompt Engineering

To address the limitations of direct predictions in complex reasoning tasks, CoT prompting emerged as a structured technique that enhances language models by guiding them through intermediate reasoning steps (Wei et al., 2022). Despite its effectiveness, CoT initially depended on manually created prompts, which constrained its adaptability. Subsequent innovations, such as Auto-CoT, addressed this by automating the creation of reasoning chains (Zhang et al., 2022). Further extensions, such as Tree-of-Thought (ToT), expanded reasoning by enabling the exploration of multiple pathways (Yao et al., 2024), while Graph-of-Thought (GoT) introduced a structured approach where reasoning steps are represented as interconnected nodes in a graph (Besta et al., 2024). These advancements directly influenced our multi-agent framework, where specialized agents leverage structured reasoning to collaboratively enhance irony detection.

2.3. Multi-agent Cooperation

Multi-agent frameworks effectively leverage collaborative interactions among specialized agents for complex tasks, with studies exploring forms like deliberation, structured debates, and dialogues. Debating frameworks, for instance, improve factual accuracy and solution diversity in complex reasoning (Du et al., 2023). Others, such as CAMEL (Li et al., 2023), use role-play to simulate nuanced human behaviors; MathChat (Wu et al., 2023b) employs structured dialogues for intricate tasks; and AutoGen (Wu et al., 2023a) shows multi-agent versatility through customizable structures. Inspired by these methodologies, our framework uses specialized LLM-based agents for contextual, semantic, and rhetorical analysis via CoT reasoning. A final decision agent integrates their outputs, enhancing irony detection’s robustness, accuracy, and interpretability.

3. Methodology

This section details our methodology for irony detection, centered on a novel multi-agent collaborative framework driven by LLM, termed the CAF-I. This framework comprises a Context Agent, a Semantic Agent, a Rhetoric Agent, a Decision Agent, and a Refinement Evaluator Agent.

CAF-I aims to address single-model limitations in irony detection via multi-agent collaboration. As irony comprehension demands nuanced understanding of context, semantics, and rhetoric, CAF-I employs dedicated agents for these dimensions. A Decision Agent aggregates these analyses, refined by the conditional feedback of the refinement evaluator agent, to achieve robust and accurate detection.

3.1. Problem Definition

Given an input text collection $X = \{x_i\}_{i=1}^n$, the primary objective is to accurately identify ironic expressions within each text x_i . This task is formalized as a binary classification problem. For each input text $x_i \in X$, the goal is to predict a Boolean label $y_i \in \{\text{Ironic}, \text{Non-Ironic}\}$, indicating the presence or absence of irony. This work aims to develop an advanced LLM-powered system to address this challenge.

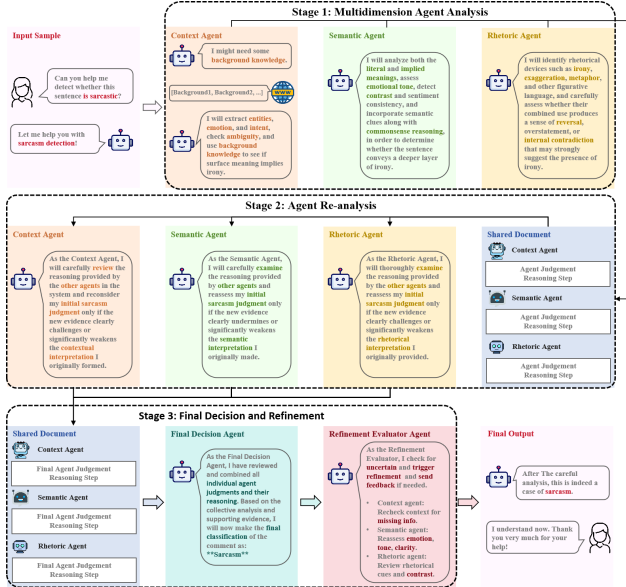


Figure 2. The Overall architecture of CAF-I. The framework includes three specialized agents, Context, Semantic, and Rhetoric, which conduct multi-round analysis and optimization of the input text. The Decision Agent integrates the preliminary classification, and the Evaluator Agent determines if a round of feedback optimization is needed.

3.2. Overall Workflow

The CAF-I inference workflow (Figure 2) unfolds in key stages: Context (CA), Semantic (SA), and Rhetoric (RA)

agents first independently provide initial assessments and reasoning. These are shared for collaborative reanalysis where agents refine judgments using peer insights. A Decision Agent (DA) then aggregates these refined judgments for an initial irony classification. Finally, a Refinement Evaluator Agent (REAgent) reviews this, potentially triggering a single conditional refinement loop with targeted feedback to improve the final classification within one iteration. The detailed design is presented in 3.3.

3.3. Agent Design

3.3.1. CONTEXT AGENT

(1) Role The Context Agent identifies irony via semantic analysis of input x_i and strategic external knowledge use, detecting inconsistencies between literal and contextual meanings.

(2) Workflow Upon receiving an input text x_i , the Context Agent extracts internal contextual features using an LLM. This involves identifying key named entities $E = \{e_1, \dots, e_k\}$, their semantic relationships $R = \{r_1, \dots, r_l\}$, and the overarching theme T . These components collectively define the internal context, denoted as $\mathcal{C}_{int} = (E, R, T)$. The overall extraction process can be formulated as:

$$\mathcal{C}_{int} = \text{ExtractInternalContextLLM}(x_i) \quad (1)$$

Subsequently, the CA, guided by an LLM, assesses the need for external knowledge based on the ambiguity of x_i and the novelty of \mathcal{C}_{int} . If deemed beneficial, a search query q is formulated using $q = g_{\text{LLM}}(\text{keywords}(x_i))$. Retrieval via an external search API yields a set of documents $D = \{d_1, \dots, d_p\}$. These are summarized into an external context representation S .

Finally, a specialized classifier, $\text{IronyClassifier}_{\text{LLM}}$, integrates the original input text x_i , the extracted internal context \mathcal{C}_{int} , and the summarized external knowledge S (if retrieved) to directly output a binary decision $y_1 \in \{\text{Ironic}, \text{Non-Ironic}\}$ regarding the presence of irony, along with associated reasoning. This process can be represented as:

$$(y_1, R_1) = \text{IronyClassifier}_{\text{LLM}}(x_i, E, R, T, S) \quad (2)$$

3.3.2. SEMANTIC AGENT

(1) Role The Semantic Agent employs a CoT process for detailed semantic analysis to identify irony. It focuses on the interplay of emotional tone, implied meaning, and commonsense expectations to detect nuanced incongruities.

(2) Workflow The workflow initiates with an LLM performing semantic parsing of the input x_i to discern its literal meaning $M_{lit}(x_i)$ from pragmatically implied intent

$I_{imp}(x_i)$:

$$(M_{lit}(x_i), I_{imp}(x_i)) = \text{SemanticParse}_{LLM}(x_i) \quad (3)$$

Subsequently, as part of its CoT reasoning, the agent identifies the expressed emotion $E_{exp}(x_i)$:

$$E_{exp}(x_i) = \text{EmotionClassifier}_{LLM}(x_i) \quad (4)$$

It then evaluates the consistency between $E_{exp}(x_i)$ and the contextually expected emotion $E_{ctx}(x_i)$. Concurrently, the SA assesses the text’s consistency with commonsense knowledge \mathcal{K}_{cs} :

$$\begin{aligned} \text{Consistency}(x_i, \mathcal{K}_{cs}) &= \text{Verifier}_{LLM}(x_i, \mathcal{K}_{cs}) \\ &\in \{\text{Consistent}, \text{Inconsistent}\}. \end{aligned} \quad (5)$$

The SA’s final irony classification y_2 and its supporting explanation R_2 are generated by an LLM that holistically reasons over the outcomes of semantic parsing, the evaluation of emotional consistency, and commonsense reasoning through a CoT process:

$$(y_2, R_2) = \text{SemanticIronyDetector}_{LLM}(M_{lit}(x_i), I_{imp}(x_i), E_{exp}(x_i), E_{ctx}(x_i), \text{Consistency}(x_i, \mathcal{K}_{cs})) \quad (6)$$

3.3.3. RHETORIC AGENT

(1) Role The Rhetoric Agent identifies irony by analyzing rhetorical devices within the input text x_i . It evaluates their interplay with communicative intent, using CoT prompting, to assess irony based on linguistic form.

(2) Workflow Employing an LLM with CoT, the RA first identifies rhetorical devices D_r present in the text:

$$D_r = \text{RhetoricalDeviceDetector}_{LLM}(x_i) = \{d_{r1}, \dots, d_{rm}\} \quad (7)$$

For each detected device $d_{rj} \in D_r$, the LLM explains its rhetorical function E_{rj} :

$$E_{rj} = \text{ExplainRhetoricalFunction}_{LLM}(d_{rj}) \quad (8)$$

Subsequently, the RA then assesses if the overall rhetorical structure, guided by D_r and their explanations $\{E_{rj}\}_{j=1}^m$, indicates irony, potentially evidenced by contradictions, reversals, or exaggerations. An LLM-based classifier generates the final judgment y_3 and explanation R_3 from these inputs:

$$(y_3, R_3) = \text{RhetoricalIronyClassifier}_{LLM}(D_r, \{E_{rj}\}_{j=1}^m) \quad (9)$$

3.3.4. DECISION AGENT

(1) Role The final Decision Agent produces a robust and interpretable irony classification y by consolidating the analysis agents’ second-round outputs, comprising their judgments and reasoning traces.

(2) Workflow The Decision Agent initiates its process by collecting the second-round outputs from the three primary analysis agents: their binary judgments $y_1^{(2)}, y_2^{(2)}, y_3^{(2)}$ and their corresponding reasoning traces $R_1^{(2)}, R_2^{(2)}, R_3^{(2)}$.

The DA employs a hierarchical strategy. If agents achieve consensus, their unanimous judgment, denoted as y^* , is adopted:

$$y = y^*, \quad \text{if } y_1^{(2)} = y_2^{(2)} = y_3^{(2)} \quad (10)$$

If a two-agent majority exists, the decision is made by majority vote:

$$\begin{aligned} y &= \text{majority}(y_1^{(2)}, y_2^{(2)}, y_3^{(2)}), \quad \text{if } (y_1^{(2)} = y_2^{(2)} \neq y_3^{(2)}) \\ &\quad \vee (y_1^{(2)} = y_3^{(2)} \neq y_2^{(2)}) \vee (y_2^{(2)} = y_3^{(2)} \neq y_1^{(2)}) \end{aligned} \quad (11)$$

In cases of complete disagreement, an LLM analyzes the agents’ reasoning traces ($R_i^{(2)}$) for clarity, coherence, and relevance, adopting the judgment supported by the most compelling argument.

The DA ultimately outputs the final binary classification $y \in \{\text{Ironic}, \text{Non-Ironic}\}$ with a synthesized justification derived from the most persuasive reasoning.

3.3.5. REFINEMENT EVALUATOR AGENT

(1) Role. The Refinement Evaluator Agent operates subsequently, once per inference, to assess the preliminary decision s and explanation \mathcal{E} from the Decision Agent. It determines if a single refinement iteration is necessary (R_{needed}) and generates hypothetical feedback as a triplet (f_{CA}, f_{SA}, f_{RA}) for upstream agents to improve reliability.

(2) Workflow. The REAgent first evaluates s and \mathcal{E} for internal quality, inferring a confidence level $C_{RE} \in \{\text{High}, \text{Medium}, \text{Low}\}$, and checks for strong contradictions, yielding $F_{contra} \in \{\text{Yes}, \text{No}\}$. The necessity of refinement, R_{needed} , is determined by:

$$R_{needed} = \begin{cases} \text{"true"} & \text{if } C_{RE} \in \{\text{Low}\} \vee F_{contra} = \text{Yes} \\ \text{"false"} & \text{otherwise} \end{cases} \quad (12)$$

If $R_{needed} = \text{"true"}$, the REAgent generates concise textual feedback components f_{CA}, f_{SA}, f_{RA} , suggesting attention points based on weaknesses inferred from s and \mathcal{E} . The REAgent outputs R_{needed} and the feedback triplet, signaling the control logic to either accept the initial result or perform a single, final refinement iteration.

4. EXPERIMENT

In this section, we try to answer the following research questions:

- **RQ1:** How does CAF-I perform on benchmark datasets against diverse existing irony detection methods?
- **RQ2:** What are the performance contributions of CAF-I’s core components?
- **RQ3:** How robust is CAF-I’s architecture with different LLM backbones, and how does it compare to standard prompting of these LLM?
- **RQ4:** How interpretable is CAF-I, and how valuable are its agents’ intermediate reasoning steps for detection performance?
- **RQ5:** How does CAF-I’s inference efficiency compare to other competitive CoT-based LLM approaches?

4.1. Experiment Setup

To evaluate the efficacy and generalizability of CAF-I, we utilized four established sarcasm detection benchmarks: IAC-V1, IAC-V2, MUSTARD, and SemEval-2018 Task 3. Model performance was assessed using Accuracy and the Macro-F1 score, with Macro-F1 serving as the primary evaluation metric. The GPT-4o model was employed as the LLM backbone for all agents within the CAF-I framework, and a temperature setting of 0 was used for all API interactions to ensure reproducibility. We compared CAF-I against a comprehensive suite of baselines, including LLM-based approaches, fine-tuned Pre-trained Language Models, and deep learning method. For complete details regarding the experimental setup, including specific dataset characteristics, a detailed description of evaluation metrics, implementation specifics, and a comprehensive list of comparison baselines, please refer to Appendix A.

4.2. Overall Performance Comparison

The overall performance comparison, with results detailed in table 1, reveals several key insights.

Our Framework Achieves New SOTA. Our proposed CAF-I framework demonstrates clear superiority, establishing a new state-of-the-art. It achieves the highest average Accuracy of 76.89 percent and an average Macro-F1 of 76.31 percent across all benchmarks, consistently outperforming other methods on most datasets and metrics. This underscores the efficacy of CAF-I’s structured multi-agent collaborative reasoning for irony detection.

CoT in CAF-I Outperforms Simpler Prompting. Within LLM-based approaches, our CAF-I framework clearly demonstrates the significant advantage of its CoT reasoning over simpler zero-shot prompting. Specifically, CAF-I’s multi-agent CoT process achieves an average Macro-F1 of 76.31 percent, substantially outperforming the GPT-4o zero-shot baseline’s average Macro-F1 of 67.75 percent. This

highlights the substantial benefits derived from employing a sophisticated CoT methodology like that inherent in CAF-I.

LLM-based Methods Outperform Non-LLM Approaches. Finally, LLM-based methods, on the whole, exhibit stronger performance compared to non-LLM approaches. Advanced LLM systems like CAF-I, with an average Macro-F1 of 76.31 percent, substantially outperform leading non-LLM methods; for example, the fine-tuned RoBERTa achieved an average Macro-F1 of 70.42 percent, and the traditional deep learning model DC-Net reached 71.33 percent. This indicates the enhanced capability of sophisticated LLM architectures for complex irony understanding over earlier paradigms.

4.3. Ablation Study

To assess the contribution of each core component within CAF-I, we conducted ablation studies, with detailed results presented in Figure 3. These studies involved systematically removing specialized analysis agents CA, SA, or RA from all analysis stages, or deactivating the REAgent’s conditional refinement mechanism, thereby making the initial aggregated decision final. Experiments were performed on the IAC-V1, MuSTARD, and SemEval-2018 datasets.

The results affirm the integral role of every component. Removing any single analysis agent consistently led to significant performance degradation across all datasets, underscoring their critical contributions. For instance, ablating the RA typically incurred a substantial average drop in Macro-F1. Similarly, deactivating the REAgent’s refinement mechanism also noticeably reduced performance, particularly on more challenging datasets like SemEval-2018, confirming the value of this step for enhancing decision robustness. These findings validate the necessity of each specialized agent and the refinement process within CAF-I.

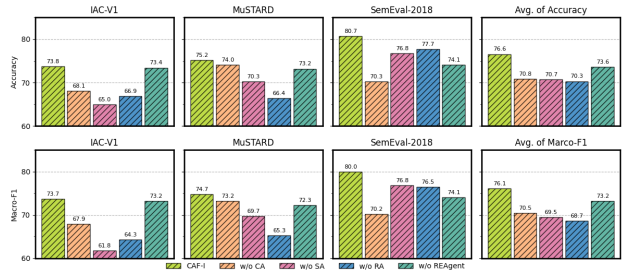


Figure 3. The ablation experiment result of CA, SA, RA and REAgent.

4.4. Robust Experiment

To assess CAF-I’s architectural robustness, we evaluated its performance using alternative LLM backbones, specifically Qwen 2-7B and Llama 3-8B, on the four benchmark datasets while maintaining the core framework and prompts. These

Table 1. Overall performance comparison across four benchmark datasets. All LLM strategies are zero-shot. Acc. denotes Accuracy and Ma-F1 signifies Macro-F1. Best results are presented in **bold**, second-best are underlined. Scores are reported as %.

Method	IAC-V1		IAC-V2		MuSTARD		SemEval-2018		Avg.	
	Acc.	Ma-F1	Acc.	Ma-F1	Acc.	Ma-F1	Acc.	Ma-F1	Acc.	Ma-F1
MIARN (Tay et al., 2018)	63.21	63.18	72.75	72.75	64.60	63.90	68.50	67.80	67.26	66.91
SAWS (Pan et al., 2020)	66.13	65.60	76.20	76.20	69.71	70.95	69.90	68.90	70.48	70.41
DC-Net (Liu et al., 2021)	66.50	66.40	78.00	77.90	<u>71.28</u>	<u>71.43</u>	70.80	69.60	<u>71.64</u>	<u>71.33</u>
BERT (Devlin et al., 2019)	65.30	65.20	76.40	76.20	64.30	64.30	69.90	68.40	68.97	68.52
RoBERTa (Liu et al., 2019)	70.10	69.90	76.60	76.70	66.10	66.00	70.20	69.10	70.75	70.42
GPT-4o (Yao et al., 2024)	70.63	70.05	73.03	71.99	67.24	65.79	64.03	63.17	68.73	67.75
GPT-4o+CoC (Yao et al., 2024)	<u>72.19</u>	<u>71.52</u>	73.36	72.31	69.42	68.48	70.79	70.60	71.44	70.73
GPT-4o+Goc (Yao et al., 2024)	65.00	62.91	64.97	61.30	70.69	69.91	<u>74.03</u>	<u>74.02</u>	68.67	67.04
GPT-4o+Boc (Yao et al., 2024)	68.75	67.36	71.35	69.39	69.42	68.45	62.12	61.85	67.91	66.76
CAF-I (Ours)	73.75	73.71	<u>77.87</u>	<u>76.82</u>	75.21	74.73	80.73	79.99	76.89	76.31

results were compared against each LLM’s baseline IO prompt performance, sourced from SarcasmCue (Yao et al., 2024), to isolate architectural contributions. The evaluation confirmed CAF-I’s architectural value. Across different LLM backbones, CAF-I consistently demonstrated substantial average Macro-F1 improvements over their respective IO prompting baselines. For instance, CAF-I with Qwen 2-7B achieved a notable 81.25% Macro-F1 on SemEval-2018, surpassing even its primary GPT-4o backbone implementation 79.99% Ma-F1 on this dataset. This general trend of outperformance strongly confirms CAF-I’s robustness and broad applicability. Detailed comparative results are presented in Appendix C.

4.5. Explainability Experiment

4.5.1. CASE STUDY

To illustrate the interpretability and error correction capabilities of CAF-I, a case study was conducted. In this instance (details are provided in Figure 4), CAF-I initially an incorrect “Sarcasm” judgment arose because the Semantic and Rhetoric Agents over-interpreted certain phrases, despite the Context Agent finding no significant inconsistency. However, the REAgent identified this discrepancy and triggered a conditional refinement. In the subsequent feedback-guided re-evaluation, the Rhetoric Agents revised their assessments, aligning with the Context Agent towards “Not Sarcastic.” This revised consensus allowed the Decision Agent to correctly classify the input, demonstrating REAgent’s effectiveness in error correction and enhancing judgment robustness.

4.5.2. QUANTITATIVE EVALUATION OF EXPLANATIONS

Furthermore, to quantitatively assess the utility of the intermediate reasoning generated by CAF-I’s specialized agents, an auxiliary experiment is conducted. This evaluation, detailed in Appendix D, confirmed that the agent-generated explanations provide valuable discriminative information, thereby enhancing the framework’s interpretability.

4.6. Inference Efficiency Analysis

To assess practical trade-offs, we analyzed CAF-I’s inference efficiency and detection performance against standard CoT and ToT baselines, all using the gpt-4o-mini backbone for fair comparison. This evaluation, detailed in Table 2, used a 400-sample subset from our four benchmarks.

Table 2 indicates CAF-I achieves the highest accuracy and Macro-F1 scores. Its average inference time of 9.67 seconds per sample is highly competitive, nearly identical to that of ToT at 9.64 seconds and only marginally more than standard CoT at 8.93 seconds, despite CAF-I’s sophisticated multi-agent architecture.

These results demonstrate CAF-I not only significantly outperforms both baselines in detection accuracy but does so without a substantial increase in computational overhead. This highlights the efficiency of CAF-I’s collaborative reasoning design, leveraging architectural depth for superior accuracy while maintaining practical inference speed.

Table 2. The average inference time (s/sample) and corresponding detection performance compared to baselines, using gpt-4o-mini on a 400-sample subset.

Method	Acc.	Ma-F1	Avg. Inf Time
CoT	63.50	62.48	8.93
ToT	68.75	66.44	9.64
CAF-I	72.50	71.03	9.67

5. CONCLUSION

This paper introduced **CAF-I**, an LLM-driven multi-agent framework for robust irony detection, integrating specialized context, semantic, and rhetoric agents with collaborative refinement. Experiments confirmed CAF-I’s state-of-the-art zero-shot performance over existing methods. Ablation studies validated its component necessity and architectural robustness, while also highlighting its interpretability and error correction.

References

- Baziotis, C., Athanasiou, N., Papalampidi, P., Kolovou, A., Paraskevopoulos, G., Ellinas, N., and Potamianos, A. Ntua-slp at semeval-2018 task 3: Tracking ironic tweets using ensembles of word and character level attentive rnns. *arXiv preprint arXiv:1804.06659*, 2018.
- Besta, M., Blach, N., Kubicek, A., Gerstenberger, R., Podstawski, M., Gianinazzi, L., Gajda, J., Lehmann, T., Niewiadomski, H., Nyczyk, P., et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 17682–17690, 2024.
- Booth, W. C. *A rhetoric of irony*. University of Chicago Press, 1974.
- Castro, S., Hazarika, D., Pérez-Rosas, V., Zimmermann, R., Mihalcea, R., and Poria, S. Towards multimodal sarcasm detection (an _obviously_ perfect paper). *arXiv preprint arXiv:1906.01815*, 2019.
- Davidov, D., Tsur, O., and Rappoport, A. Semi-supervised recognition of sarcasm in twitter and amazon. In *Proceedings of the fourteenth conference on computational natural language learning*, pp. 107–116, 2010.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.
- Du, Y., Li, S., Torralba, A., Tenenbaum, J. B., and Mordatch, I. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2023.
- Ghosh, A. and Veale, T. Fracking sarcasm using neural network. In *Proceedings of the 7th workshop on computational approaches to subjectivity, sentiment and social media analysis*, pp. 161–169, 2016.
- Kader, F. B., Nujat, N. H., Sogir, T. B., Kabir, M., Mahmud, H., and Hasan, K. Computational sarcasm analysis on social media: a systematic review. *arXiv preprint arXiv:2209.06170*, 2022.
- Li, G., Hammoud, H., Itani, H., Khizbullin, D., and Ghanem, B. Camel: Communicative agents for” mind” exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008, 2023.
- Liang, B., Lou, C., Li, X., Yang, M., Gui, L., He, Y., Pei, W., and Xu, R. Multi-modal sarcasm detection via cross-modal graph convolutional network. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1767–1777, 2022.
- Liu, X., Yu, H., Zhang, H., Xu, Y., Lei, X., Lai, H., Gu, Y., Ding, H., Men, K., Yang, K., et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Liu, Y., Wang, Y., Sun, A., Meng, X., Li, J., and Guo, J. A dual-channel framework for sarcasm recognition by detecting sentiment conflict. *arXiv preprint arXiv:2109.03587*, 2021.
- Lukin, S. and Walker, M. Really? well. apparently bootstrapping improves the performance of sarcasm and nastiness classifiers for online dialogue. *NAACL 2013*, pp. 30, 2013.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- Oraby, S., Harrison, V., Reed, L., Hernandez, E., Riloff, E., and Walker, M. Creating and characterizing a diverse corpus of sarcasm in dialogue. *arXiv preprint arXiv:1709.05404*, 2017.
- Pan, H., Lin, Z., Fu, P., and Wang, W. Modeling the incongruity between sentence snippets for sarcasm detection. In *ECAI 2020*, pp. 2132–2139. IOS Press, 2020.
- Pennington, J., Socher, R., and Manning, C. D. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543, 2014.
- Poria, S., Cambria, E., Hazarika, D., and Vij, P. A deeper look into sarcastic tweets using deep convolutional neural networks. *arXiv preprint arXiv:1610.08815*, 2016.
- Qiu, Z., Yu, J., Zhang, Y., Lai, H., Rao, Y., Su, Q., and Yin, J. Detecting emotional incongruity of sarcasm by commonsense reasoning. *arXiv preprint arXiv:2412.12808*, 2024.
- Reyes, A., Rosso, P., and Veale, T. A multidimensional approach for detecting irony in twitter. *Language resources and evaluation*, 47:239–268, 2013.
- Tay, Y., Tuan, L. A., Hui, S. C., and Su, J. Reasoning with sarcasm by reading in-between. *arXiv preprint arXiv:1805.02856*, 2018.

- Van Hee, C., Lefever, E., and Hoste, V. Semeval-2018 task 3: Irony detection in english tweets. In *Proceedings of the 12th international workshop on semantic evaluation*, pp. 39–50, 2018.
- Wallace, B. C., Charniak, E., et al. Sparse, contextually informed models for irony detection: Exploiting user communities, entities and sentiment. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1035–1044, 2015.
- Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q. V., Zhou, D., et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Wu, Q., Bansal, G., Zhang, J., Wu, Y., Li, B., Zhu, E., Jiang, L., Zhang, X., Zhang, S., Liu, J., et al. Autogen: Enabling next-gen llm applications via multi-agent conversation. *arXiv preprint arXiv:2308.08155*, 2023a.
- Wu, Y., Jia, F., Zhang, S., Li, H., Zhu, E., Wang, Y., Lee, Y. T., Peng, R., Wu, Q., and Wang, C. Mathchat: Converse to tackle challenging math problems with llm agents. *arXiv preprint arXiv:2306.01337*, 2023b.
- Xue, H., Xu, L., Tong, Y., Li, R., Lin, J., and Jiang, D. Breakthrough from nuance and inconsistency: Enhancing multimodal sarcasm detection with context-aware self-attention fusion and word weight calculation. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 2493–2503, 2024.
- Yao, B., Zhang, Y., Li, Q., and Qin, J. Is sarcasm detection a step-by-step reasoning process in large language models? *arXiv preprint arXiv:2407.12725*, 2024.
- Zhang, M., Zhang, Y., and Fu, G. Tweet sarcasm detection using deep neural network. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: technical papers*, pp. 2449–2460, 2016.
- Zhang, Y., Zou, C., Lian, Z., Tiwari, P., and Qin, J. Sarcasm-bench: Towards evaluating large language models on sarcasm understanding. *arXiv preprint arXiv:2408.11319*, 2024.
- Zhang, Z., Zhang, A., Li, M., and Smola, A. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.

A. Detail Experiment Set Up

A.1. Datasets

We evaluate CAF-I’s efficacy and generalizability using four established sarcasm detection benchmarks:

- **IAC-V1** (Lukin & Walker, 2013): This dataset comprises comments from online political debates, offering examples of sarcasm embedded within argumentative discourse.
- **IAC-V2** (Oraby et al., 2017): An extension of IAC-V1, providing a larger, more diverse collection of sarcastic/non-sarcastic statements from similar online discussion forums.
- **MuSTARD** (Castro et al., 2019): Sourced from popular television shows, MuSTARD provides conversational context for each utterance.
- **SemEval-2018 Task 3** (Van Hee et al., 2018): Consists of English tweets annotated for irony, representing challenges of social media text like brevity, informal language, and implicit contextual cues for ironic intent.

Table 3 summarizes detailed dataset statistics. Their diversity in domains and challenges enables robust framework evaluation.

Table 3. Overview of the benchmark datasets used for evaluating irony detection.

Dataset	Year	Size	Avg. Length	Domain	Context
IAC-V1	2013	320	68	Debate	No
IAC-V2	2016	1042	43	Debate	No
MuSTARD	2019	784	14	Dialogue	Yes
SemEval-2018	2018	183	14	Twitter	No

A.2. Evaluation Metrics

We evaluate model performance using **Accuracy** for overall correctness and the **Macro-F1 score** as the primary metric, following standard practices in the field (Yao et al., 2024). Macro-F1 averages the F1-scores for ironic and non-ironic classes independently, providing a balanced measure suitable for potentially imbalanced datasets common in irony detection.

A.3. Experiment Details

The **GPT-4o** model serves as the LLM backbone for all agents within the proposed CAF-I framework, accessed via the official OpenAI API. Its selection ensures state-of-the-art baseline capabilities in natural language understanding and reasoning. To enhance reproducibility, a temperature setting of 0 was used for all API interactions.

A.4. Comparison Baselines

To contextualize CAF-I’s performance, we compare it against a comprehensive suite of baselines across three categories: **LLM-based:** Approaches from the SarcasmCue framework (Yao et al., 2024), including GPT-4o Zero-Shot and three advanced prompting strategies: Chain of Contradiction, Graph of Cues, and Bagging of Cues. Performance figures are cited from their reported results. **Fine-tuned PLMs:** Standard Pre-trained Language Models such as BERT-base-uncased (Devlin et al., 2019) and RoBERTa-base (Liu et al., 2019), fine-tuned on target datasets. **Deep Learning Methods:** Influential deep learning methods with task-specific architectures, namely MIARN (Tay et al., 2018), SAWS (Pan et al., 2020), and DC-Net (Liu et al., 2021). Results are sourced from literature with aligned evaluation settings (Qiu et al., 2024; Zhang et al., 2024; Tay et al., 2018; Xue et al., 2024; Pan et al., 2020).

B. Case Study Illustration

This appendix section presents a detailed visual walkthrough of the case study discussed in Section 4.5.1. Figure 4 below illustrates the step-by-step multi-agent reasoning, collaborative reanalysis, and conditional refinement process undertaken by the CAF-I framework on an example input. This visualization aims to further clarify CAF-I’s interpretability and its mechanism for error correction, as highlighted in the main text, by showing how the agents interact and revise their judgments to reach a final, more accurate classification.

Input: Actually there are already a couple floating around. One fascinating one: Gaia theory and the Daisyworld model. It's all around self-organising systems and the earth as a self-regulating system/organism/something or other. James Lovelock and co. Interesting to read but sometimes feels a little too 'lets-all-just-get-along' to be real, hard-core theory. Either way, it's probably a lot more likely to be accepted than creationism, should evolutionism for some reason be disproved. Label: **Not Sarcasm**	
Round One	Round Two
Context Agent Stage1: Not sarcasm. No significant inconsistency was identified between the statement and its context. The text discusses the Gaia Theory and Daisyworld model in an exploratory manner, noting its self-organizing concepts and contrasting its acceptance with creationism. While the description includes subjective judgments, there is no implied meaning that directly contradicts the surface meaning. Additionally, the context is insufficient to infer sarcasm confidently. ✔	Context Agent Stage1: Not sarcasm. No significant inconsistency or irony was identified between the statement and its context. The discussion of Gaia theory and its scientific positioning—including the comparison to creationism—appears genuine rather than sarcastic. The hypothetical tone and critique reflect critical engagement rather than implied mockery or contradiction. The context is sufficient to reject sarcasm. ✔
Semantic Agent Stage1: Is sarcasm. The statement employs a playful or skeptical tone, especially with phrases like "feels a little too 'lets-all-just-get-along' to be real, hard-core theory" and the comparison of Gaia theory to creationism. There is a subtle contradiction or mismatch between the seemingly serious discussion of theories and the implied dismissal of Gaia theory as overly idealistic or unscientific, suggesting an inversion of meaning typical of sarcasm. ✘	Semantic Agent Stage1: Is sarcasm. The statement exhibits a subtle sarcasm due to the sharp mismatch between the implied dismissiveness of the Gaia theory ("feels a little too 'lets-all-just-get-along' to be real, hard-core theory") and the seemingly neutral information it presents. The tone suggests skepticism about the seriousness of the theory while pretending to give it credit, creating an inversion of meaning typical of sarcasm. ✘
Rhetorical Agent Stage1: Is sarcasm. The phrase "sometimes feels a little too 'lets-all-just-get-along' to be real, hard-core theory" employs **irony** and **hyperbole** , mocking the Gaia theory's perceived idealism by contrasting its supposed depth or seriousness with a trivialized and exaggerated tone of naivety. This contradiction creates a clear rhetorical cue for sarcasm. ✘	Rhetorical Agent Stage1: Not sarcasm. The statement provides a detailed explanation of Gaia theory and contrasts its potential acceptance with creationism without employing clear rhetorical devices like irony, hyperbole, or tonal inversion that would signal sarcasm. While the phrase "feels a little too 'lets-all-just-get-along' to be real, hard-core theory" might indicate mild skepticism or critique, it is literal and not mockingly contradictory; thus, the rhetoric supports a non-sarcastic reading. ✔
Context Agent Stage2: Not sarcasm. My initial Context Analysis reasoning remains valid despite reviewing shared inputs. The shared reasoning from the Rhetorical Agent emphasizes hyperbole and irony, but this evidence does not create a fundamental contradiction or gap within the surface-level meaning of the statement. While the tone includes subjective judgments, it does not exhibit a clear inversion of intent or a direct clash with the literal meaning necessary for confident sarcasm detection from a contextual perspective. ✔	Context Agent Stage2: Not sarcasm. My initial Context Analysis reasoning remains valid despite reviewing shared inputs. Neither Semantic Analysis nor Rhetorical Analysis provides evidence of sarcasm that contradicts my earlier perspective. The critique of the Gaia theory's tone ("feels a little too 'lets-all-just-get-along'...") reflects genuine skepticism rather than a sarcastic inversion of meaning, and the broader context supports a literal, non-sarcastic tone. ✔
Semantic Agent Stage2: Is sarcasm. My initial Semantic Analysis reasoning remains valid after reviewing the shared inputs. The phrase "feels a little too 'lets-all-just-get-along' to be real, hard-core theory" uses hyperbole to subtly undermine the seriousness of Gaia theory. This inversion of tone aligns with sarcastic cues, and the Rhetorical Agent's analysis reinforces my original judgment without revealing any fundamental contradiction to invalidate it. ✘	Semantic Agent Stage2: Is sarcasm. My initial reasoning remains valid despite reviewing shared inputs. The phrase "feels a little too 'lets-all-just-get-along' to be real, hard-core theory" exhibits a tonal incongruity and dismissiveness that aligns with semantic markers for sarcasm. The shared inputs lack sufficient evidence to overturn this assessment, as they interpret the tone too literally and overlook the subtle inversion present in the statement. ✘
Rhetorical Agent Stage2: Is sarcasm. My initial Rhetorical Analysis reasoning remains valid despite reviewing shared inputs. The phrase "feels a little too 'lets-all-just-get-along' to be real, hard-core theory" still demonstrates **irony** by downplaying the Gaia theory's idealism with a tone of playful skepticism. This tonal inversion aligns with rhetorical sarcasm techniques, and no new evidence in the shared reasoning fundamentally contradicts this interpretation. ✘	Rhetorical Agent Stage2: Not sarcasm. My initial Rhetorical Analysis reasoning remains valid. While the Semantic Analysis suggests sarcasm due to tonal inversion, this perspective does not align with the rhetorical devices necessary to confirm sarcasm. The critique of Gaia theory is literal, and there is no clear use of irony, hyperbole, or mockingly contradictory language indicative of sarcasm in the statement. ✔
Final Decision Agent: Is sarcasm. Semantic analysis identified a contradiction in tone, which was explicitly supported by Rhetorical analysis highlighting irony and tonal inversion. While Contextual analysis rejected sarcasm, its reasoning lacked specific refutation of the credible evidence provided by the other two agents. ✘	Final Decision Agent: Not sarcasm. Chose NO based on lack of credible and corroborated evidence. Semantic analysis suggesting sarcasm (tonal incongruity) was uncorroborated by Contextual and Rhetorical analyses, which both firmly interpreted the critique as literal and non-sarcastic, with no compelling counter-evidence to refute their positions. ✔
Refinement Evaluator Agent: Need Reflection. Context agent needs to refine contextual analysis to explicitly address irony-related evidence and provide a more specific refutation if sarcasm is rejected. Rhetorical Agent needs to Clarify how tonal inversion and irony align clearly with the decision to avoid ambiguities in semantic contradiction claims.	

Figure 4. Illustration of the CAF-I framework’s step-by-step reasoning and refinement process on a input example.

C. Detailed Robustness Evaluation with Different LLM Backbones

This section provides the detailed performance metrics for the CAF-I framework when implemented with alternative Large Language Model backbones, namely Qwen 2-7B and Llama 3-8B. The results, as summarized in Table 4, are compared against each LLM’s baseline Input-Output prompt performance. This data supports the discussion in Section 4.4 regarding CAF-I’s architectural robustness and its ability to enhance performance across different underlying models.

Table 4. Robustness evaluation of the CAF-I framework using different LLM backbones Qwen 2-7B, Llama 3-8B compared to the IO prompting. Baseline IO results are aligned with SarcasmCue (Yao et al., 2024). All scores are in %.

LLM Backbone	IAC-V1		IAC-V2		SemEval-2018		MuSTARD	
	Acc.	Ma-F1	Acc.	Ma-F1	Acc.	Ma-F1	Acc.	Ma-F1
Qwen2-7b + IO	56.56	49.32	51.85	38.57	45.15	38.83	54.78	46.17
Qwen2-7b + CAF-I	71.85	71.19	67.08	66.39	81.63	81.25	72.60	73.80
Llama3-8b + IO	55.94	46.40	54.70	43.74	49.36	44.46	54.64	44.99
Llama3-8b + CAF-I	60.62	60.40	70.36	73.55	75.76	75.47	59.56	57.11

D. Quantitative Evaluation of Agent Explanations

An auxiliary experiment on the SemEval-2018 dataset, using GPT-3.5 Turbo as a baseline, assessed the utility of intermediate reasoning from CAF-I’s analysis agents. We compared GPT-3.5 Turbo’s performance under two conditions: using only

the original input text with a standard IO prompt, versus augmenting the input with a feature prompt concatenating stage 1 textual explanations from CAF-I’s Context, Semantic, and Rhetoric agents. Incorporating these agent explanations improved the baseline model’s **Macro-F1** score from **68.4%** to **70.2%**. This performance gain indicates that CAF-I’s specialized agents generate valuable, discriminative reasoning, substantiating the framework’s effective internal processing and enhancing its interpretability.

E. Future Work

Future research could focus on extending CAF-I’s application beyond irony detection to include other forms of emotional detection and broader nuanced linguistic phenomena, as well as adapting it for cross-lingual applications. Further work might also involve developing the framework for multimodal contexts, enhancing agent collaboration dynamics, and systematically evaluating its performance and efficiency with a wider array of Large Language Model backbones.