



TONE TRACER: DECODING THE HIDDEN MEANING

¹Anshika Mishra, ²Anurag Shrivastava, ³Ishaan Tripathi, ⁴Riddhisha Srivastava, ⁵Sachin Tripathi

¹B. Tech Student, ²Professor, ³B. Tech Student, ⁴B. Tech Student, ⁵B. Tech Student

¹Engineering (All Branches),

¹Babu Banarsi Das Northern India Institute of Technology, ¹Lucknow, ¹India

Abstract: Sarcasm detection is tricky because sarcastic expressions usually carry meanings that are the opposite of their literal words. This often leads to traditional sentiment analysis systems missing the mark. This work presents Tone Tracer, a sarcasm detection framework that considers context and uses transformer-based language models. The method looks at both an utterance and its conversational context to identify implicit contradictions and changes in tone. Experimental results on benchmark datasets show better precision, recall, and F1-score compared to traditional and context agnostic methods. The study emphasizes the need for context when it comes to detecting sarcasm accurately.

Index Terms— Sarcasm Detection, Context-Aware NLP, Transformer Models, Sentiment Analysis, Conversational Context.

I. INTRODUCTION

Sarcasm detection is a tough task in natural language processing because of the gap between literal meaning and the intended feeling. This issue impacts sentiment analysis systems. Traditional rule-based and feature driven methods are not very reliable [1][2]. Many deep learning approaches also overlook the context of conversations. While transformer models enhance semantic representation, looking at sentences alone is not enough to grasp sarcastic intent. To tackle this problem, we present Tone Tracer. This framework uses a context aware transformer-based method to encode both target utterances and conversational context together. It shows better results on standard sarcasm detection datasets.

II. RELATED WORK

Sarcasm detection has received a lot of attention in NLP because it affects sentiment analysis and understanding figurative language. Previous studies have explored various methods, including rule-based, machine learning, deep learning, transformer based, context aware, and multimodal approaches. Each method targets different aspects of the problem.

2.1 Text-Based and Traditional Approaches

Early methods relied on rule-based and lexicon-driven techniques [1]. They used signals like sentiment incongruity, punctuation, and intensifiers, but struggled to generalize. Later, traditional machine learning models employed handcrafted features and distant supervision. These models improved scalability but were still limited by shallow semantics and noisy labels.

2.2 Deep Learning Approaches for Sarcasm Detection

Deep learning models such as LSTM, BiLSTM, and CNN decreased dependence on manual features by learning distributed representations [3]. However, most of these models processed utterances in isolation, which limited their ability to capture sarcasm that relies on context at the discourse level [13].

2.3 Transformer Based Models

Model based on transformers like BERT and RoBERTa enhanced sarcasm detection by capturing rich contextual and semantic relationships [2][4]. Despite their strong performance, many still treat sarcasm detection as a sentence-level task and overlook the conversational context.

2.4 Context Aware Sarcasm Detection

Context aware methods use conversational history or discourse information [3][8][11] to better understand implicit sarcastic intent. Jointly encoding context and target utterances with transformers has led to consistent improvements. However, dealing with irrelevant or noisy context remains a challenge.

2.5 Multimodal Sarcasm Detection

Multimodal approaches blend text with visual or audio signals [10] to model cross-modal incongruity, especially in social media settings. Although these methods are effective, they often require complex architectures and large annotated datasets, which limits scalability and generalization.

2.6 Research Gap

Current methods either overlook contextual cues, lack robustness, or become overly complex with multimodal inputs. This highlights the need for a scalable, text-based, context aware transformer framework for effective sarcasm detection.

III. TASK DEFINITION

Sarcasm detection is treated as a supervised binary text classification task. The goal is to predict if a specific statement shows sarcastic intent. In a conversation that consists of a series of statements, the model learns to connect the target statement with its earlier context to assign a sarcasm label. Unlike sentence-level sentiment classification, sarcasm detection needs to capture hidden practical meaning and semantic inconsistencies that can contradict the obvious sentiment. This work uses a context aware approach to model the relationships in discourse that are important for correctly identifying sarcasm.

IV. DATASET DESCRIPTION

Experiments are carried out on sarcasm detection datasets that consist of user-generated text from online platforms. These datasets include both sarcastic and non-sarcastic examples. They are annotated using manual and distant supervision methods.

4.1 Data Sources

Data is gathered from social media and online discussion platforms like Twitter and forums [3][8], where sarcastic language frequently appears. Each example contains a short text with optional preceding conversational context.

4.2 Annotation Strategy

Sarcasm labels are obtained using distant supervision with sarcasm-related hashtags and through manual annotation. Distant supervision allows for large-scale data collection but adds noise, while manual annotation improves quality but can be subjective.

4.3 Context Construction

Conversational context is created by combining a fixed number of utterances that come before the target text. Only prior context is used to avoid information leakage. Sentence-level inputs are used when context is not available.

4.4 Dataset Statistics

The datasets show class imbalance, with non-sarcastic examples occurring more often than sarcastic ones. The data is divided into training, validation, and test sets using stratified sampling in Table 1.

Table 1 Dataset Statistics

Feature	Description	Value
Data Source	Twitter & Online Discussion Forums	Combined
Sarcastic Samples	Label 1 (Positive Class)	62,813
Non-Sarcastic Samples	Label 0 (Negative Class)	18,595
Total Samples	Total records used for study	81,408
Annotation Method	Distant Supervision & Manual	Hybrid

4.5 Preprocessing

Text is pre-processed by lowercasing, removing URLs and mentions, and normalizing whitespace. Tokenization and context concatenation are done using the pretrained transformer tokenizer with special separator tokens.

V. PROPOSED METHODOLOGY

This section introduces Tone Tracer, a framework for detecting sarcasm that considers context and is based on transformer models. The method models a target utterance along with its conversational context using a complete architecture without any manually designed features.

5.1 System Overview

Tone Tracer includes input construction, transformer-based encoding, and a classification layer. The model estimates whether a target utterance is sarcastic by encoding it together with the preceding context.

5.2 Input Representation

The conversational context and target utterance are combined into a single input sequence with special separator tokens. This setup allows the model to understand interactions between the context and the target text through self-attention.

5.3 Context Aware Encoding Using Transformers

The combined input is processed by a pretrained transformer encoder [6][7][14], following standard fine-tuning practices for text classification [15], to create contextualized token embeddings.

5.4 Sarcasm Classification Layer

The [CLS] embedding goes through a feed-forward layer followed by SoftMax to predict sarcasm probabilities. This layer connects contextual representations to sarcastic and non-sarcastic categories.

5.5 Training Objective

The model is trained using cross-entropy loss to reduce prediction errors. Class weighted loss is used to address imbalances in sarcasm datasets.

5.6 Optimization Strategy

Training uses the AdamW optimizer [6] with a learning rate scheduler. Early stopping and gradient clipping help maintain stable and efficient convergence.

5.7 Context Agnostic Baseline Configuration

A baseline model is created using only the target utterance without any context. Comparing the performance of this model with the context aware version shows the value of conversational context.

VI. EXPERIMENTAL SETUP

This section outlines the experimental setup used to test the Tone Tracer framework. All experiments take place in controlled settings to ensure a fair comparison with baseline models.

6.1 Data Splits

Each dataset is divided into training, validation, and test sets using stratified sampling to keep class distribution intact. The training set is for optimization, the validation set is for tuning and early stopping, and the test set is for final evaluation only in Table 2.

Table 2 Dataset Split Distribution

Split Type	Distribution (%)	Number of Samples	Primary Purpose
Training Set	80%	40,000	Model optimization and weight updates
Validation Set	10%	5,000	Hyperparameter tuning and early stopping
Test Set	10%	5,000	Final performance and generalization check
Total	100%	50,000	-

6.2 Baseline Models

Tone Tracer is compared with traditional machine learning models, deep learning classifiers, and transformer-based sentence-level baselines. These models help assess how contextual modelling improves performance across different architectures.

6.3 Hyperparameter Settings

Transformer based models are fine-tuned with a consistent set of hyperparameters chosen based on validation performance. All tuning occurs on the validation set, and dropout is used to reduce overfitting.

6.4 Training Configuration

Models are trained with mini-batch gradient descent, using early stopping based on validation loss. Class weighted loss and several random seeds address class imbalance and enhance result stability.

6.5 Implementation Details

All models are built in Python using PyTorch and the Hugging Face Transformers library. Tokenization is done with the corresponding pretrained tokenizer, and training runs on GPU-enabled hardware.

6.6 Reproducibility

Reproducibility is maintained through fixed random seeds, consistent preprocessing, and standardized evaluation methods. Model checkpoints are saved based on validation performance throughout all experiments.

VII. EVALUATION METRICS

Model performance is assessed using standard classification metrics. Precision, recall, and F1-score are the main measures because of class imbalance in sarcasm datasets. Precision shows how accurate the predicted sarcastic instances are. Recall indicates the model's ability to recognize actual sarcasm. The macro-averaged F1-score serves as the primary evaluation metric to ensure balanced performance across classes, while accuracy is provided as an additional metric. This evaluation protocol follows common practices in sarcasm detection research.

VIII. RESULTS AND ANALYSIS

This section presents the experimental results of the Tone Tracer framework and examines its performance. The context aware transformer model is compared to traditional, deep learning, and context agnostic transformer baselines.

8.1 Quantitative Results

Table 3 summarizes model performance based on accuracy, precision, recall, and F1-score. The context aware model consistently outperforms all baselines. The largest gains are seen in F1-score and recall for sarcastic instances.

Table 3 Quantitative Performance Comparison of Sarcasm Detection Models

Metric Category	Accuracy	Precision	Recall	F1-Score
Overall Performance	71.0%	0.66	0.69	0.66
Regular (Class 0)	-	0.86	0.72	0.79

Metric Category	Accuracy	Precision	Recall	F1-Score
Sarcastic (Class 1)	-	0.45	0.67	0.54

8.2 Comparison with Baseline Models

Traditional machine learning models perform poorly because they rely on handcrafted features and shallow representations. Deep learning and sentence-level transformer baselines show better results but still fall short compared to the proposed approach. This outcome emphasizes the importance of explicit contextual modelling.

8.3 Impact of Contextual Modelling

The context aware model achieves higher F1-scores than its context agnostic counterpart, confirming the value of conversational context. Improvements are most noticeable for implicit sarcasm that lacks clear lexical indicators.

8.4 Error Analysis

Error analysis reveals problems in situations that require knowledge of the outside world or cultural references. Very subtle sarcasm and literal statements with exaggerated phrasing also result in false negatives and false positives.

8.5 Discussion of Results

Overall, the results show that sarcasm detection greatly benefits from context aware transformer modelling. The findings support the idea that sarcasm is a practical phenomenon that needs discourse-level understanding instead of just focusing on individual sentences.

IX. DISCUSSION

The results show that detecting sarcasm relies heavily on the conversation's context. Sarcastic intent usually comes from differences in meaning or use, not just single statements. Using models that consider context greatly improves the ability to identify sarcastic remarks. This helps catch subtle sarcasm that does not have clear language cues. Comparing these models with those that ignore context proves that understanding the meaning of sentences alone doesn't work for detecting sarcasm in conversation. An analysis of errors also points out problems linked to knowledge about the world and overlapping figures of speech. This suggests a need to include more detailed contextual and external information in future work.

X. LIMITATIONS

Despite its improved performance, the proposed Tone Tracer framework has several limitations. The model relies only on textual input and cannot capture multimodal cues like visual or acoustic signals, which often convey sarcasm in real-world settings. Its effectiveness is also limited by noisy and subjective annotations as well as the availability and quality of conversational context. Furthermore, using transformer -based architectures adds higher computational demands, which restricts its use in resource-limited or real-time environments.

XI. APPLICATIONS

Detecting sarcasm is important for NLP applications that need good sentiment and intent understanding. The proposed Tone Tracer framework can improve sentiment analysis and opinion mining by cutting down on misclassifications in user-generated content. In social media monitoring and conversational AI systems, handling sarcasm enables better opinion tracking and more suitable responses. Moreover, sarcasm detection helps with content moderation by finding harmful or passive-aggressive language that might get past regular filters.

XII. FUTURE WORK

Future work can build on this framework by adding multimodal signals like images or audio to better capture sarcasm that comes from mixed signals. Another important area is improving generalization across different domains through domain adaptation and few-shot learning approaches [12]. Including external knowledge and commonsense reasoning may also help detect sarcasm that depends on implicit assumptions. Lastly, it is crucial to improve model efficiency and clarity for real-time use and clear decision-making.

XIII. CONCLUSION

This paper introduced Tone Tracer, a framework for sarcasm detection that is aware of context and uses a transformer-based approach. The experimental results show that considering context significantly improves the detection of subtle sarcasm, which sentence-level models often misclassify. The findings support the idea that sarcasm is a practical phenomenon that occurs at the discourse level, not just at the word level. Although there are still challenges to address, Tone Tracer provides a solid foundation for future work on multimodal, knowledge-aware, and efficient systems for detecting sarcasm.

References

- [1] Bhattacharyya, P., Kumar, A., & Banerjee, S. (2017). *A survey on sarcasm detection*. ACM Computing Surveys, 50(5), 1–38.
- [2] Joshi, A., Sharma, V., & Bhattacharyya, P. (2016). *Harnessing context incongruity for sarcasm detection*. ACL 2016, 757–762.
- [3] Dong, Y., Chen, J., & Yang, X. (2020). *Modeling conversational context for sarcasm detection*. AAAI Conference on Artificial Intelligence, 748–755.
- [4] Ghosh, D., & Veale, T. (2016). *Fracking sarcasm using neural network*. Computational Linguistics, 42(2), 297–327.
- [5] Mishra, S., Dandapat, P., & Ekbal, A. (2020). *Figurative language identification using contextualized word representations*. Figurative Language Processing Workshop, 97–102.
- [6] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). *BERT: Pre-training of deep bidirectional transformers for language understanding*. NAACL-HLT, 4171–4186.
- [7] Liu, Y., Ott, M., Goyal, N., et al. (2019). *RoBERTa: A robustly optimized BERT pretraining approach*. arXiv preprint arXiv:1907.11692.
- [8] Zhang, R., & Zhu, H. (2019). *Sarcasm detection with context-aware multi-task learning*. WWW Conference, 3581–3587.
- [9] Majumder, N., et al. (2019). *Dialogue RNN: An attentive RNN for emotion detection in conversations*. AAAI, 6818–6825.
- [10] Poria, S., et al. (2019). *Context-dependent sentiment analysis in user-generated videos*. ACL, 873–883.
- [11] Hazarika, D., et al. (2018). *Modeling inter-utterance dependency for sarcasm detection*. NAACL, 281–286.
- [12] Schick, T., & Schütze, H. (2021). *Exploiting cloze questions for few-shot text classification*. EACL, 255–269.
- [13] Tang, D., Qin, B., & Liu, T. (2015). *Document modeling with gated recurrent neural network for sentiment classification*. EMNLP, 1422–1432.
- [14] Vaswani, A., et al. (2017). *Attention is all you need*. NeurIPS, 5998–6008.
- [15] Sun, C., Qiu, X., Xu, Y., & Huang, X. (2019). *How to fine-tune BERT for text classification?* China National Conference on Chinese Computational Linguistics, 194–206.