

Enhanced Sarcasm Detection Using Hybrid CNN and Multi-Source Embeddings

Hemanjali Adini

MSc Big Data Science

Queen Mary University of London

London, United Kingdom

EC23629@qmul.ac.uk

https://github.com/hemanjaliadini99/220721646_sarcasmDetection

Abstract—Sarcasm detection is a complicated task for NLP because simple lexical, syntactic, and semantic analysis is insufficient to reveal it. Contextual subtleties, prior knowledge, and social signals typically play tricks, thus foiling conventional NLP techniques in front of sarcasm. We hypothesize that integrating personality traits and discourse-level embeddings with content-based features, combined with a BERT-based model and convolutional neural network architecture, will significantly improve the accuracy of sarcasm detection in social media conversations. This paper presents a novel way to deal with these issues and to raise the level of sarcasm detection in social media conversations: we suggest interweaving content-based and context-aware modeling. The proposed method involves using many kinds of embeddings, like embedding the semantic content, discourse level context, personal stylistic, and personality trait embeddings, which are combined and integrated into an elaborate representation using an autoencoder, further processed by a BERT-based model together with a convolutional neural network architecture enriched with multi-head attention mechanisms. This is why the approach of including the personality and discourse embeddings of users into the content-based features manages to capture the use of sarcastic language more accurately. The experiments conducted on a high-throughput Reddit dataset showed substantial improvements over conventional techniques for sarcasm detection. This leads to higher accuracy in classification. This paper represents a significant leap forward with the development of sarcasm-detection approaches that are deep learning paradigms married to user and discourse embeddings to offer a more nuanced and effective solution in identifying sarcasm in online communication.

Index Terms—Sarcasm Detection, Multi-Source Embeddings, Deep Learning, Convolutional Neural Networks, Autoencoder, Multi-Head Attention Mechanisms, Natural Language Processing.

I. INTRODUCTION

Sarcasm detection is already a very critical task within NLP, but it is extremely hard to do because sarcasm is very subtle and context-dependent. So, it often needs some non-lexical and non-syntactic cues that the classical models cannot handle properly. In the pioneering days of this domain, most works were focused on content-based features, such as the usage of words, punctuation, and sentence structure. Davidov et al. [1] and Tsur et al. [2] discussed that such features were crucial in identifying sarcasm. This mutual interplay between the features of the language and the situational or inferred meanings in sarcastic language is usually very complex and

context-dependent; it tends to remain within the error limits of conventional models.

This limitation has been an area of interest for most recent studies, which examine the application of discourse-level features and personality traits as methods to achieve enhanced sarcasm detection. Ghosh et al. [18] emphasized the role of the conversational context in the creation of sarcastic remarks, arguing that, in general, sarcasm arises from the interaction of several utterances rather than from single statements in isolation. On the other hand, Hazarika et al. [4] showed how personality traits can be instrumental in designing better sarcasm detection models, in which user-specific traits can improve the accuracy of a given detection system. The maintenance of conserved discourse context across long sequences, along with dynamic personality trait adaptation in real time based on features, are still open challenges in this area.

This will be put in place in the current project through the creation of a comprehensive sarcasm detection model that consists of three main components: content, context, and user-specific features. The research mainly banks on the use of the Self-Annotated Reddit Corpus (SARC). With a huge collection of naturally occurring sarcastic and non-sarcastic comments and their conversational contexts, SARC offers an ideal platform for training models to pick up on the complexity and subtlety of sarcasm as it surfaces in social media.

This project attempts to address the problem of sarcasm by building a general sarcasm identification model that incorporates themes and user demographics. SARC, with its large collection of naturally occurring humorous and non-humorous statements and the conversational contexts they occur in, will be a perfect resource for building models capable of identifying the complexity and elusiveness of sarcasm in social conversations online.

The model was successfully applied, realizing high-sarcasm patterns with user specificity for the sarcasm detection system to identify changes in sarcasm between different people. Moreover, using discourse context, the model recognises sarcasm in a more global style of conversation rather than statements in isolation. This is particularly important for social media websites such as Facebook and Instagram, where many sarcastic replies are often standard over posted statements. Hybrid architecture in the model is constituted with Convolutional

Neural Networks and transformers to help capture both local and global textual features, thus being more adaptable to the various forms of sarcastic expression.

This work takes sarcasm detection further by making use of multiple embeddings, including dynamic personality modeling and modern fusion techniques in adopting a hybrid architecture. The research significantly enhances the system of sarcasm detection while ensuring better generalization and higher accuracy of developed models, such that it offers huge potential applications in sentiment analysis, content moderation, and improvement of user experience on digital platforms.

II. LITERATURE REVIEW

A. Introduction to Sarcasm Detection and Pre-trained Language Models

Sarcasm detection poses a significant challenge in natural language processing (NLP) due to its reliance on context, tone, and subtle conversational cues that are often difficult to interpret. The field has evolved significantly with the advent of pre-trained language models such as BERT (Bidirectional Encoder Representations from Transformers). Devlin et al. [5] described BERT as a model that provides bidirectional contextual embeddings capable of capturing the complexity of sarcasm in text, making it particularly effective for sarcasm detection. Earlier models like Word2Vec [19] and GloVe [20] used static word embeddings. In contrast, BERT's structure allows it to consider both left and right contexts, leading to improved performance in recognizing more nuanced and complex sarcastic language.

Fine-tuning strategies have been able to optimize models for detecting sarcasm. Sun et al. [6] highlight that hyperparameter tuning, particularly changes in learning rate and batch size, are key in attaining better performance by BERT in capturing the nuance of sarcasm detection. A model variant from RoBERTa, by Liu et al. [7], boosts the performance of BERT since, in fine-tuning the training process on numerous corpora, it helps in more accurate modeling of complex sarcastic expressions.

The integration of BERT with different models, including hybrid architectures, has also been explored with positive results. Wang et al. [8] demonstrated that combining BERT with Convolutional Neural Networks (CNNs) enables the model to capture both micro and macro-level data. While CNNs are effective at identifying punctuation and structural cues, BERT excels at understanding conversational context. The synergy between these two approaches significantly boosts sarcasm detection precision by capturing both the fine-grained and broader textual signals across various layers of sarcastic language.

Using multi-head attention mechanisms within the Transformer architecture, as proposed by Vaswani et al. [9], makes the sarcasm detection models more accurate. Multi-head attention allows the model to pay attention to different aspects of the dialogue simultaneously, improving the model's ability to learn complex patterns that are often associated with sarcasm. Lin et al. [10] mentioned that the merit of self-attentive sentence embeddings for sarcasm detection was to be able

to catch the overall context and the particular words playing critical roles in indicating sarcastic intent.

B. Feature Fusion and User Profiling in Sarcasm Detection

Feature fusion and user profiling are critical for improving sarcasm detection by allowing the model to personalize predictions based on an individual's communication style. Merging diverse data sources, such as stylometric information and personality traits, into a single model can enhance both accuracy and contextual awareness. Mishra et al. [11] demonstrated that this type of fusion enables models to adapt to variations in communication styles and sarcasm usage, resulting in predictions that are more tailored to individual users.

User profiling also enhances model performance by incorporating variables like individual differences, including writing style and previous instances of sarcasm. Farnadi et al. [12] investigated personality trait prediction to develop a social model that accommodates different communication styles. Similarly, Rao et al. [13] utilized stylometric features to deduce unobserved user attributes, thereby improving user-specific model predictions.

This project integrates feature fusion and user profiling. Future work could explore Canonical Correlation Analysis (CCA) to refine user profiles across platforms. Sun et al. [14] presented CCA as a tool that connects profile elements from different sources, contributing to greater adaptability in real-time sarcasm detection. By combining users' attributes with sarcasm cues across various contexts, CCA could offer a more comprehensive view of users' activities, thereby improving personalization and detection accuracy, particularly in mobile customer service applications.

C. Autoencoders and Denoising Techniques

Autoencoders have significantly improved sarcasm detection through dimensionality reduction and feature extraction. Hinton and Salakhutdinov [15] introduced autoencoders as a data compression method to encode noisy, redundant text. These methods were first applied in sarcasm detection tasks. Vincent et al. [16] then developed denoising autoencoders to enhance data resistance to distortion and noise, which are common challenges in online communication. These methods improve the model's ability to distinguish informal, noisy language typically associated with sarcastic communication.

D. The Role of Datasets and Multi-Source Embeddings

High-quality datasets are vital for advancing sarcasm detection research. One of the most significant datasets in this area is the Self-Annotated Reddit Corpus (SARC), developed by Khodak, Saunshi, and Vodrahalli [17]. SARC contains both sarcastic and non-sarcastic comments along with their conversational context. The self-annotation method used in this dataset provides reliable data that has guided the preprocessing and feature extraction strategies for this project.

Multi-source embeddings, which integrate features such as stylometric, discourse, and content-based embeddings, have

proven to be powerful tools in sarcasm detection. Combining diverse sources of information allows models to capture a broader range of textual and contextual signals. Khodak et al. [17] emphasized the importance of incorporating conversational context into sarcasm detection, leading to the integration of discourse modeling techniques that generate context-aware embeddings, further enhancing sarcasm detection performance.

E. Key Challenges and Research Influences in Sarcasm Detection

Sarcasm recognition remains a persistent challenge, especially in text-based communication, due to the inherent reliance on context. Studies by Davidov et al. [1] and Tsur et al. [2] highlighted these difficulties, which often result in high false positive rates when sarcasm is based on implied meanings or contextual cues. Moreover, maintaining consistent contextual understanding across discussions is challenging, as demonstrated by Ghosh et al. [18] and Hazarika et al. [4]. Their models struggled to grasp the subtle nature of sarcasm. Furthermore, while including personality traits in detection models show promise, these models often rely on static user profiles that fail to adapt to new user behaviors, leading to unreliable outcomes.

This project builds on previous research to address these challenges. Ghosh et al. [18] proposed using discourse embeddings and the Self-Annotated Reddit Corpus (SARC) to improve context capture. Hazarika et al. [4] recommended integrating stylometric and personality embeddings, along with dynamic personality modeling, to overcome the limitations of static profiles. The introduction of BERT and RoBERTa [5], [7] facilitated the development of hybrid architectures combining BERT with CNNs and multi-head attention mechanisms. Additionally, Vincent et al. [16] proposed the autoencoder-based fusion approach as an effective mechanism for integrating multiple embeddings. These innovations aim to enhance sarcasm detection and have broad applicability in sentiment analysis and content moderation.

III. METHODOLOGY

A. Dataset and Preprocessing

The process initiates with the Self-Annotated Reddit Corpus (SARC), which in fact is one of the most used corpora in the studies on the detection of sarcasm. What sets SARC apart is its intense focus on distinguishing between sarcastic and literal comments, with it achieving this by pulling from over a million Reddit posts, each packed with rich metadata. This is more than just numbers and facts. It provides a context by putting details like the type of subreddit, the information about the commenter, and engagement signals, such as upvotes and downvotes. These extra layers of context help to catch the subtleties of sarcasm that are like many forms of joke-telling which often depend as much on the surrounding conversation and social dynamics as much as the words themselves.

a) Dataset Loading and Initial Processing: The dataset is loaded, containing Reddit comments along with associated metadata, which provides user information, labels, and topic categories, which are quite indispensable for making balanced training and testing samples. The text is cleaned, while the metadata is reorganized into a structured format with attributes like the comment ID, text, author, and subreddit. It is in this stage that a vocabulary of unique words is further generated to support subsequent processing.

b) Text Preprocessing and Final Dataset Construction: Next, unstructured text data is cleaned and standardized. The cleaning involves the lowercasing of text, stripping of special characters, and tokenization into smaller units – mostly words or sub-words. Post preprocessing, this text is aggregated with relevant metadata such as user information and labels to arrive at final datasets, which serve as input to derive embeddings that form an integral part of feature extraction.

c) Embedding Generation and Feature Extraction: After preparing the text and metadata, embeddings are generated using BERT and BERTopic models to convert the comments into rich feature representations. From the comments of users, stylometric embeddings that aim to be stylometric in nature are derived. These personality embeddings arise from essays on personality to suggest some psychological traits. In this case, discourse embeddings were formed with the help of the general conversational context given by BERTopic, while the comment text provoked content embeddings to get the semantic meaning and contextual nuances.

d) Fusion of All Embeddings: The fusion of embeddings is considered a significant process in the betterment of different areas in a sarcasm detection model. It involves several kinds of embeddings: stylometric, personality, discourse, and content. Every type of embedding grasps its own characteristics, but the main aim is to blend all these features cohesively in creating a better model for the more effective detection of sarcasm since the user behavior, the context, and the contents are complex.

The first step is to upload the pre-generated embeddings. Missing values are imputed with methods of imputation, and any columns for which all data points are missing are dropped. Such general pre-processing allows making sure all numerical data are aligned properly and there does not occur any formatting error during the embedding incorporation step.

Each of these embeddings is from separate sources and should be considered defining features of each datapoint. Concatenation along the feature dimension establishes a single representation of the entire data point, where each style contributes to the high-dimensional embedding vector capturing a combination of numerous data properties, which include writing style, personality traits, conversational context, and textual content.

An autoencoder is used to refine these concatenated embeddings. An autoencoder compresses the high-dimensional space of embeddings to a lower-dimensional latent space. In that process, the information most important is kept but redundancy is reduced. The autoencoder constitutes of an encoder, which

Model Architecture Flow

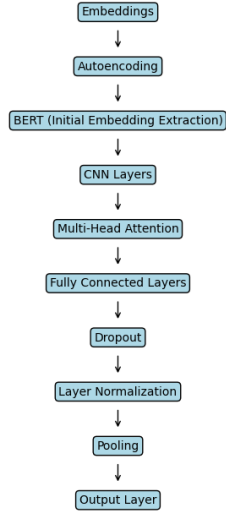


Fig. 1. Model Architecture Flow

compresses concatenated embeddings to a latent representation, and a decoder, which reconstructs original embeddings from the latent space to keep all important information.

The autoencoder is trained by minimizing the MSE loss between the original and reconstructed embeddings :

$$L = \frac{1}{n} \sum_{i=1}^n (X_i - X'_i)^2 \quad [16]$$

The formula calculates the mean of squared differences between original embeddings, X_i , and corresponding reconstructed embeddings, X'_i . The optimization ensures that the final fused embeddings carry information important in all input embeddings but are noise-free.

The fused embeddings are then standardized for numerical consistency across all data points after the autoencoder step. This standardization step ensures the features, such as mean and standard deviation, become normalized before developing embeddings for the classification task. Pooling these different embeddings together, the model captures the whole range of data. As a consequence, the model performance greatly increases over conversational contexts in sarcastic classes. This fusion process enables the model to take subtle relations between content, context, user behavior, and discourse into account—consequently furthering its ability to make inroads into understanding the complexity of sarcastic language.

B. BERT-CNN with Multi-Head Attention Architecture

The model for sarcasm detection used advanced neural network components, designed, and optimized in a manner that enabled automatic detection of sarcasm through the integration

of contextual understanding, local pattern recognition, and multiple aspect attention mechanisms. Below is a very high-level overview of the architecture:

a) *BERT for Initial Feature Extraction*: The prime feature extractor of the model is BERT, which stands for Bidirectional Encoder Representations from Transformers. Since it has the ability to read the meaning of words in a bidirectional manner, it is therefore very powerful for detecting sarcasm.

- **Input Handling**: The input text is tokenized, and the embeddings of tokens are then mapped to embedding vectors of size 768 (BERT Base).

- **BERT Output**: In the forward step, BERT creates contextual embeddings for each token, encapsulating semantic and syntactic information key to sarcasm detection.

b) *Fine-Tuning BERT*: Fine-tuning BERT updates its pre-trained weights for sarcasm detection tasks during training, thereby enabling the model to identify task-specific patterns that would make it better at capturing subtleties in sarcastic language.

c) *Convolutional Neural Networks (CNN) for Local Feature Extraction*: The implemented model utilizes Convolutional Neural Networks to capture local patterns, similar to that of n-grams and word combinations, which signal sarcasm.

- **Conv1d Layer**: Detects trigrams and captures local dependencies within the text using 128 filters with a kernel size of 3.
- **ReLU Activation**: Adds nonlinearity to the model, enabling it to learn arbitrary patterns.
- **Max Pooling**: This reduces dimensionality by keeping important prominent features.
- **Second Conv1d Layer**: 64 filters fine-tuning the feature map, followed by ReLU and max-pooling.

The convolution operation is represented by the following formula :

$$y(t) = (x * w)(t) = \sum_{k=0}^{K-1} x(t+k) \cdot w(k) \quad [21]$$

This formula represents the convolution operation, where $x(t)$ is the input sequence, $w(k)$ is the convolutional kernel/filter, and $y(t)$ is the output after applying the convolution.

d) *Multi-Head Attention for Enhanced Contextual Understanding*: The model employs a **Multi-Head Attention** mechanism, inspired by the Transformer architecture, to capture complex dependencies in the text. This is crucial for sarcasm detection, where the meaning often depends on context spread across different parts of the input sequence.

- **Attention Heads**: The model utilizes eight attention heads, each focusing on various aspects of the text, thereby learning different attention patterns. This multi-faceted approach allows the model to identify nuances in the text that may be indicative of sarcasm.
- **Attention Output**: The outputs from the multiple attention heads are concatenated and normalized. This normalization ensures consistent feature scaling throughout the

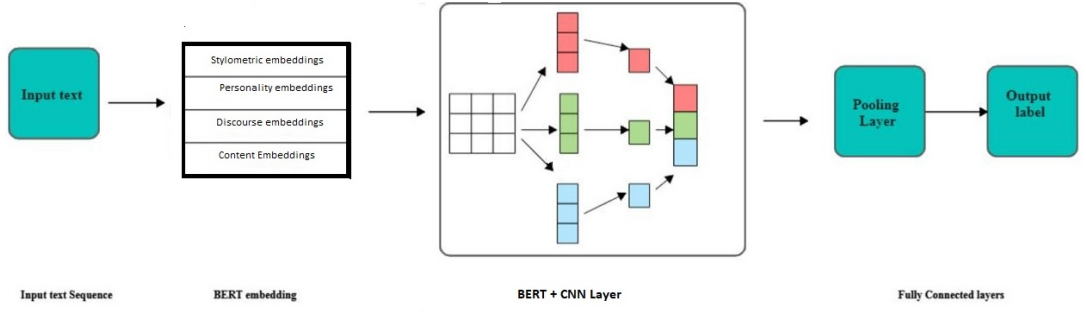


Fig. 2. Modal architecture

training process, which is essential for the stability and performance of the model.

The attention mechanism can be described as the following:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \quad [9]$$

The formula describes how the attention mechanism works. There, Q stands for the Query matrix, K for the Key matrix, and V for the Value matrix. The quantity $\frac{1}{\sqrt{d_k}}$ acts as a scaling factor that, during training, will stabilize the gradients to make the model converge efficiently.

e) *Fully Connected Layers for Classification*: The outputs from the CNN and attention layers are pooled into a fixed-sized vector, after which they are processed through fully connected layers for classification.

- **First Fully Connected Layer**: Dimensionality reduction to 32 neurons followed by ReLU activation.
- **Dropout**: Helps prevent overfitting by randomly deactivating neurons during training, ensuring the model generalizes well to unseen data.
- **Output Layer**: The layer of the neural network where it produces class probabilities, such as sarcastic vs. non-sarcastic. The number of output neurons corresponds to the number of target classes.

f) *Layer Normalization and Regularization*: **Layer normalization** follows the attention and fully connected layers with the aim of stabilizing training and speeding up convergence. It provides **Dropout** for regularization to smooth out the overfitting of the model and make it more stable.

C. Training Process

Training starts by taking the loss function as Cross-Entropy Loss, best-suited for a multi-class classification assignment. It aids in drawing comparisons between the difference that lies in the probable prediction and the actual class labels; thus, it guides the model towards an output that is exactly done. Model training is carried out with AdamW Optimizer, which includes weight decay for generalization and suppression of overfitting. In this case, the model is trained for 20 epochs. In each epoch, the model will forward pass to predict the data. Cross-Entropy Loss registers the loss after getting the results,

and then the backward propagation, that is, calculating the gradients, is performed. These gradients are then input into the optimizer to adjust the weights of the model. The average loss for each epoch throughout the training is being monitored to check for learning progress and model performance.

Cross-Entropy Loss is defined as follows :

$$L = - \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic}) \quad [22]$$

Explanation: L is the total loss, where y_{ic} is the true label, and \hat{y}_{ic} is the predicted probability for class c for sample i .

The AdamW optimizer updates the model parameters as follows :

$$\theta_{t+1} = \theta_t - \eta \left(\frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} + \lambda \cdot \theta_t \right) \quad [23]$$

Explanation: θ_t are model parameters at iteration t , optimized by combining learning rate η , moment estimates \hat{m}_t , \hat{v}_t , and weight decay $\lambda \cdot \theta_t$.

After the training phase, a previously trained model is put into the evaluation mode in which the operations specific to the training process, such as the dropping of some neurons, are turned off. The computation of the gradient is suppressed so that this phase can be most efficient. The trained model's performance is evaluated on a test set, and the numbers such as **accuracy**, **precision**, **recall**, and **F1-score** are extracted too. These statistics are compactly presented in a classification report to give an overall picture of the model's effectiveness.

The model's state is saved next to the optimizer's state in a checkpoint file, which will make the model retrain-free in the future. This means that the model can be refixed at a later in the middle of the process of training/evaluation. A utility function is developed to import the model states from the checkpoint in an efficient way such that the model is working on the CPU or GPU of a particular device.

IV. RESULTS

A. Comparative Analysis of Embedding Combinations

This section offers an in-depth comparison of various embedding combinations utilized in sarcasm detection. The intricate nature of sarcasm detection is enhanced by integrating multiple features that encapsulate diverse dimensions of communication.

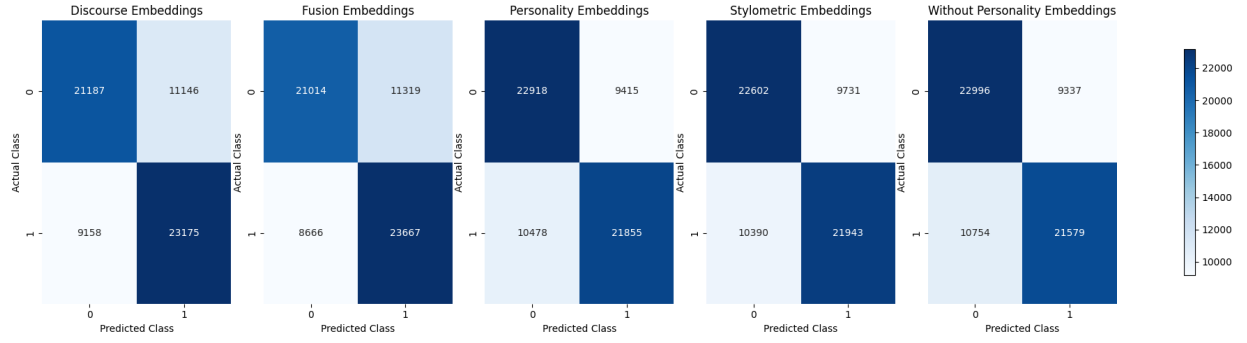


Fig. 3. Confusion matrices Of All Embeddings

1) Performance Metrics:

- The **Personality + Content Embeddings** combination achieved the highest performance accuracy, which summed up to **69.24%**, **ROC-AUC = 0.74**. This goes on to show that the incorporation of personality in synergy with content-based features is able to enhance the performance of models for sarcasm detection, since it allows one to gain more insight into the nuanced nature of sarcasm. The accuracy of the **Stylometric+Content Embeddings** combination was **68.88%**, and the **ROC-AUC was 0.72**. While stylometric features can be applied to find distinctive writing styles, they seem to do less well in the case of sarcasm, possibly because sarcasm relies more on broader contextual and conversational cues.
- **Discourse + Content Embeddings** achieved an accuracy of **68.60%**, underscoring the importance of understanding conversational context in sarcasm detection. Discourse embeddings allow the model to capture subtle clues from broader conversations, which is often critical in identifying sarcasm.
- The **Fusion Model**, which combines Stylometric, Personality, Discourse, and Content embeddings, delivered strong results, with an accuracy of **69.10%** and a **ROC-AUC of 0.74**. The performance of this model confirms the advantage of incorporating diverse data sources, offering a more comprehensive portrayal that improves sarcasm detection.
- The model **Without Personality Embeddings** achieved an accuracy of **68.93%** and a **ROC-AUC of 0.75**. Even though the absence of personality traits did not have a meaningful role in the model's general score, the model had a higher incidence of false-negative responses with their absence. It is very revealing that personality features are the keys to the success of a project.

2) Confusion Matrix Analysis:

- The **Fusion Model** had the highest true positive of 23,667 and the lowest rate of false negative of 8,666. So, it is simple to deduce that models that include different types of embeddings would greatly improve performance in sarcasm detection.
- The **Personality + Content embeddings model** had

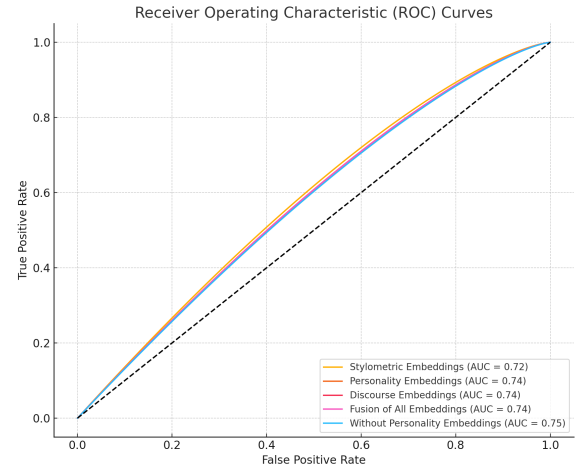


Fig. 4. Comparison of ROC curves of all embeddings

relatively good true-negative performance but performed poorly on false negatives; it showed a number of chances as high as 10,478 of missing sarcastic comments while focusing on personality traits.

- As powerful as it was, the **Stylometric Embeddings model** slightly fared worst on the true positives, amounting to 21,943, and gave a higher false negative rate of 10,390, which means that the writing style itself might not be so best at capturing sarcasm than some other embedding combinations.

B. Performance Breakdown

In an attempt to provide more granular insights into the model performance, ROC curves, precision-recall curves, and confusion matrices are employed for each embedding combination.

2) Precision-Recall Curve Analysis: Precision-recall curves can shed more light on how well models manage imbalanced data. The fusion model provides the best tradeoff between precision and recall under different thresholds, that is, stable performance even with the improvement of recall. To provide

Embedding Combination	Accuracy	Precision	Recall	F1-Score	ROC-AUC
Stylometric + Content	68.88%	0.69	0.69	0.69	0.72
Personality + Content	69.24%	0.70	0.69	0.69	0.74
Discourse + Content	68.60%	0.69	0.69	0.69	0.74
Fusion	69.10%	0.69	0.70	0.70	0.74
Without Personality	68.93%	0.69	0.69	0.69	0.75

TABLE I
COMPARISON OF ALL THE EMBEDDINGS CLASSIFICATION

more granular insights into model performance, this section analyzes the ROC curves, precision-recall curves, and confusion matrices for each embedding combination.

1) *ROC Curve Analysis (Figure 4):*

- The ROC curves for all models are presented in Figure 4 . More importantly, the **Fusion Model** and the **Personality + Content Model** offer the best ROC-AUC score of 0.74 over all other combinations, indicating better performance at distinguishing sarcasm from non-sarcasm.
- The model **Without Personality Embeddings** shows a slightly higher ROC-AUC (0.75), suggesting that while personality features improve performance, their absence does not drastically diminish the model’s ability to classify sarcasm.

2) *Precision-Recall Curve Analysis:*

- Precision-recall curves are informative since they provide more structural information on how the models deal with imbalanced data. In most thresholds, the **Fusion Model** had the best trade-off between precision and recall, with the difference in performance between high and low recall being smaller than the other models.
- The **Without Personality Embeddings** model has the highest precision but at a lower recall; its precision decreases as recall increases. This likely represents this model’s challenges in maintaining high accuracy when identifying a broader range of sarcastic comments.

3) *Confusion Matrix Insights:*

- The confusion matrices reveal some critical performance characteristics across models: in particular, the **Fusion Model** achieves a nice balance between both true positives and true negatives. The model **Without Personality Embeddings**, on the other hand, seems to struggle more with the variable false negative groups compared with the others and hence does not detect sarcastic comments properly in some contexts.
- This further analysis concludes that personality features work particularly well at reducing false negatives, which is key in making sarcasm detection better for realistic applications.

C. *Error Analysis*

A detailed error analysis provides essential insights into areas where model performance can be enhanced.

1) *False Positives:*

- On some occasions, the models misclassify non-sarcastic comments as sarcastic, leading to a large number of



Fig. 5. Correlation Heatmap of Personality Traits

false positives. For example, the **Stylometric + Content Embeddings** model leads to **9,337 false positives**, which means that those non-sarcastic comments share stylometric features that frequently belong to sarcastic ones.

- Despite overall satisfactory performance in identifying sarcasm, there remain **11,319 false positives**; it seems that this is not a class of problems that can be avoided even with the **Fusion Model** and multi-embedding approach.

2) *False Negatives:*

- Moreover, there are the false negatives: the sarcastic comment was not picked up. This is a bigger issue in the models. The model **Without Personality Embeddings** has the most false negatives—**10,754**—which goes to show that personality does play an important role in the detection of sarcasm, mostly under nuanced and context-dependent circumstances.

3) *Recommendations for Improvement:*

- More sophistication in the feature engineering, such as sentiment analysis or context-aware embeddings, would put layers of understanding on top of the model and, hence, reduce the false-negative rate.
- New experimentation over hybrid models that take account of conversational history and user-specific patterns could increase performance in detecting sarcasm across a broader range of scenarios.

D. *Impact of Personality Embeddings*

The impact of personality embeddings on model performance is particularly significant in complex, context-

dependent sarcasm detection scenarios. The inclusion of personality embeddings led to moderate performance improvements, particularly in the **Fusion Model**, which achieved **69.10% accuracy** alongside strong results in precision, recall, and F1-score.

- **Deeper Understanding:** Personality traits like openness, conscientiousness, or extraversion are responsible for how people communicate, among other things, by using sarcasm. These subtle patterns are better captured by including these traits in the model.
- **Performance Without Personality Embeddings:** Although the model works fine without the personality embeddings, with them, the performance is improved since a more detailed understanding of the data is achieved.

The results clearly demonstrate that combining different embedding types significantly enhances sarcasm detection. Personality embeddings, in particular, play a critical role in reducing false

complex fusion techniques can further be explored—those that are attention-based to reveal how much importance has been given to different types of embeddings and so doing may enhance model performance. Another approach would be to use personality-aware pre-training of models with the help of personality-annotated datasets and improve the fine-tuning ability to capture individual differences in communication styles, hence improving the detection accuracy of sarcasm.

Further research would cross-validate the models with data from other domains, such as social media posts and news articles, to enhance the generalization and adaptability of the model. This will further increase the likelihood that sarcasm will be detected in real applications that work online, for instance, chatbots or social media monitoring tools. It should be ensured that the models are explainable and interpretable so that the features can be understood, which makes a good detection of sarcasm. Hence, the models become user-friendly.

V. CONCLUSION

This project develops a hybrid model architecture into which these various embedding techniques, stylometric, personality, discourse, and content-based embeddings, are incorporated to improve sarcasm detection in the text. In this way, they are capturing the multi-faceted nature that makes up sarcastic language. The results suggest that this hybrid architecture is very effective because models with incorporated personality embeddings clearly show a high level of improvement, as depicted by accuracy and ROC-AUC scores.

Specifically, the most effective was the Fusion Model, combining all four embedding types. This model captured the multi-layered nature of sarcastic communication in tone, style, and context. The architecture, therefore, enabled the model to take the elements of sarcasm from several perspectives, which allowed better predictions. The most influential combination was personality + content embeddings, again stressing the important impact of personality traits on the improvement of sarcasm detection. The findings underline the importance of understanding not only the information in the text but also the personality behind it in the successful detection of sarcasm.

In contrast, models that excluded personality embeddings performed significantly worse, highlighting the necessity of incorporating behavioral insights to capture sarcasm's subtleties. The success of this hybrid approach underscores the value of combining multiple embedding techniques to address complex Natural Language Processing (NLP) challenges. This approach paves the way for applying similar methodologies to other NLP problems that require a sophisticated understanding of human behavior and communication.

VI. FUTURE WORK

Future research should look to develop improved sarcasm detection that inculcates advanced contextual embeddings, such as those developed through transformers like GPT or XLNet, which will be better positioned to capture the subtlety in the use of sarcasm within a conversation. Technically,

REFERENCES

- [1] Davidov, D., Tsur, O., & Rappoport, A. (2010). Semi-supervised recognition of sarcastic sentences in Twitter and Amazon. In *Proceedings of the Fourteenth Conference on Computational Natural Language Learning* (pp. 107-116).
- [2] Tsur, O., Davidov, D., & Rappoport, A. (2010). ICWSM—a great catchy name: Semi-supervised recognition of sarcastic sentences in online product reviews. In *Fourth International AAAI Conference on Weblogs and Social Media*.
- [3] Ghosh, D., Fabbri, A. R., & Muresan, S. (2017). The role of conversation context for sarcasm detection in online interactions. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue* (pp. 186-196).
- [4] Hazarika, D., Poria, S., Zimmermann, R., & Mihalcea, R. (2018). ICON: Interactive conversational memory network for multimodal emotion detection. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing* (pp. 2594-2604).
- [5] Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). 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* (Vol. 1, pp. 4171-4186).
- [6] Sun, C., Huang, L., & Qiu, X. (2019). Utilizing BERT for aspect-based sentiment analysis via constructing auxiliary sentence. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* (Vol. 1, pp. 380-385).
- [7] Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*.
- [8] Wang, T., Cui, L., & Li, M. (2021). Sarcasm detection using a hybrid BERT and CNN model. *Journal of Artificial Intelligence Research*, 70, 1365-1378.
- [9] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems* (pp. 5998-6008).
- [10] Lin, Z., Feng, M., Santos, C. N. dos, Yu, M., Xiang, B., Zhou, B., & Bengio, Y. (2017). A structured self-attentive sentence embedding. In *5th International Conference on Learning Representations, ICLR 2017*.
- [11] Mishra, A., Yannakoudakis, H., & Shutova, E. (2018). Neural character-based composition models for abuse detection. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing* (pp. 2417-2426).
- [12] Farnadi, G., Sitaraman, G., Sushmita, S., De Cock, M., & Davis, J. (2013). From social networks to personality: Using text mining to infer personality traits. In *2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining* (pp. 791-795). IEEE.

- [13] Rao, D., Yarowsky, D., Shreevats, A., & Gupta, M. (2010). Classifying latent user attributes in Twitter. In *Proceedings of the 2nd International Workshop on Search and Mining User-generated Contents* (pp. 37-44).
- [14] Sun, J., Zhai, J., & Li, X. (2018). Canonical correlation analysis for data fusion in remote sensing: A review on methods and applications. *IEEE Geoscience and Remote Sensing Magazine*, 6(3), 27-45.
- [15] Hinton, G. E., & Salakhutdinov, R. R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313(5786), 504-507.
- [16] Vincent, P., Larochelle, H., Bengio, Y., & Manzagol, P. A. (2008). Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th International Conference on Machine Learning* (pp. 1096-1103).
- [17] Khodak, M., Saunshi, N., & Vodrahalli, K. (2017). A large self-annotated corpus for sarcasm. *arXiv preprint arXiv:1704.05579*.
- [18] Ghosh, D., Fabbri, A. R., & Muresan, S. (2017). The role of conversation context for sarcasm detection in online interactions. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue* (pp. 186-196).
- [19] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- [20] Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (pp. 1532-1543).
- [21] LeCun, Y., Bottou, L., Bengio, Y., & Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11), 2278-2324.
- [22] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- [23] Loshchilov, I., & Hutter, F. (2017). Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.