

# SARCASM DETECTION

Manogna Vennela Ramireddy Yamini Durga Loya Velangani Joshita Lavanya Dudla

Master of Science in Data Science,  
University of New Haven,  
{mrami7, yloya1, vdudl1} @unh.newhaven.edu

## Abstract

This paper presents a study on sarcasm detection in text, employing three state-of-the-art natural language processing (NLP) models: BiLSTM, DistilBERT, and RoBERTa. The study aims to evaluate the performance of each model on diverse datasets, including the Sarcasm Headlines Dataset, Sarcasm on Reddit, and a combined dataset from Hugging Face. The models were evaluated for their ability to detect sarcasm across varied text characteristics. BiLSTM leverages GloVe embeddings and captures temporal dependencies in text, while DistilBERT offers a lightweight transformer approach, and RoBERTa focuses on robust, optimized transformer-based performance. The results were evaluated using standard metrics such as accuracy, precision, recall, and F1-score. The paper discusses the results, comparing them with baseline models, and offers an analysis of the challenges and successes in sarcasm detection. This research contributes to improving sentiment analysis, enhancing NLP applications, and advancing model architectures for challenging NLP tasks. The paper concludes with insights into the model performance, highlighting areas for further improvement and potential applications.

Code-GitHub-Link

[https://github.com/Joshita1306/sarcasm-detection\\_dsci6004\\_nlp](https://github.com/Joshita1306/sarcasm-detection_dsci6004_nlp)

## 1 Introduction

Sarcasm detection in text is a complex and critical challenge in natural language processing (NLP), as it often relies on subtle cues, context, and tone, which are difficult for computational models to accurately capture. Sarcastic remarks typically convey the opposite of their literal meaning, and their interpretation can be highly context-dependent. This makes sarcasm particularly difficult for traditional sentiment analysis techniques, which rely on direct word sentiment to classify text. With the increasing prevalence of social media platforms and online

communication, accurate sarcasm detection has become essential for various NLP applications, such as sentiment analysis, opinion mining, and conversational AI.

The goal of this study is to develop a robust sarcasm detection model by leveraging three different NLP approaches: Bidirectional Long Short-Term Memory (BiLSTM), DistilBERT, and RoBERTa. These models were chosen to evaluate different techniques for contextual understanding, sequence modeling, and transformer-based architectures. Each model offers distinct advantages in terms of handling language nuances. BiLSTM, as a recurrent neural network (RNN), captures both past and future context in sequential data, making it ideal for text with temporal dependencies. DistilBERT, a smaller and more efficient version of the BERT transformer, balances computational efficiency and performance, making it an attractive choice for large-scale applications. RoBERTa, an advanced version of BERT optimized for better performance, has been shown to outperform other transformer models in various NLP benchmarks, providing a strong foundation for sarcasm detection.

This project aims to evaluate and compare these models on three diverse sarcasm datasets: the Sarcasm Headlines Dataset, Sarcasm on Reddit, and a combined dataset from Hugging Face. The use of these varied datasets allows us to assess how well each model performs across different text types, ranging from sarcastic news headlines to casual online conversations. By analyzing these models' ability to detect sarcasm in different contexts, we will gain insights into how they generalize across domains and text structures.

In terms of methodology, the project employs several key techniques, including tokenization, GloVe embeddings for BiLSTM, and pre-trained tokenizers for DistilBERT and RoBERTa. Data preprocessing steps such as padding, normalization, and sequence standardization were applied to ensure the models

were trained on clean, consistent data. The models were evaluated using key metrics such as accuracy, precision, recall, and F1-score, with additional visual analysis through word clouds and confusion matrices to better understand model performance.

Ultimately, this research contributes to improving sentiment analysis by enhancing the ability to detect sarcasm, a critical aspect of understanding human

## 2 Related Work

Sarcasm detection in natural language processing (NLP) has been a challenging yet increasingly important research area. Understanding sarcasm is crucial for accurately interpreting sentiment in text, especially in domains like social media, customer reviews, and online forums, where sarcastic comments can significantly alter the meaning of statements. Several studies have addressed the problem of sarcasm detection, leveraging various NLP techniques ranging from traditional machine learning models to more advanced deep learning architectures.

Early work in sarcasm detection often relied on rule-based approaches or feature-based machine learning models. For example, Riloff et al. (2013) proposed a system that used a combination of lexical and syntactic features, such as negation and discourse markers, to identify sarcastic statements in Twitter posts. Their approach showed promise but was limited by the reliance on handcrafted features, which often failed to generalize across different datasets.

With the rise of deep learning, researchers began exploring neural network models for sarcasm detection. One such study by Prabhat et al. (2018) employed deep neural networks with word embeddings (such as Word2Vec) to detect sarcasm in social media text. They demonstrated that deep learning models outperformed traditional feature-based methods, highlighting the potential of neural networks to capture semantic relationships and contextual cues. However, this approach still faced challenges in handling the nuanced and complex nature of sarcasm.

In more recent work, transformer-based models such as BERT (Bidirectional Encoder Representations from Transformers) have been successfully applied to various NLP tasks, including sarcasm detection. Joshi

et al. (2020) investigated the use of BERT for sarcasm detection and showed that pre-trained language models, fine-tuned on sarcasm-labeled datasets, significantly improved performance. These models leverage large-scale contextual information, making them highly effective at capturing the intricate dependencies involved in sarcastic remarks. Additionally, RoBERTa, an optimized version of BERT, has been found to outperform BERT in several NLP tasks (Liu et al., 2019), suggesting its potential for sarcasm detection as well.

Another important contribution to sarcasm detection came from the use of recurrent neural networks (RNNs) and variants like Long Short-Term Memory (LSTM) networks. The BiLSTM (Bidirectional LSTM) model, which captures both past and future contexts in text, has shown promising results in sequence-based tasks. For example, Zhang et al. (2018) employed a BiLSTM network for sarcasm detection in social media posts, achieving competitive results compared to traditional models. BiLSTMs have an advantage in tasks like sarcasm detection, where understanding the context of surrounding words is crucial for accurately interpreting meaning.

In the domain of multi-modal sarcasm detection, some studies have also integrated visual or acoustic features to complement text-based models. For instance, Ghosh et al. (2017) combined textual and visual features to detect sarcasm in online video comments. While multimodal approaches have the potential to improve sarcasm detection, text-only models continue to dominate the field, particularly in applications where visual data is not available.

The use of pre-trained embeddings, such as GloVe (Global Vectors for Word Representation), has also been explored in sarcasm detection. GloVe embeddings capture the semantic meaning of words based on their co-occurrence in large corpora and are often used as input to deep learning models. The combination of GloVe with LSTM networks has been

The use of pre-trained embeddings, such as GloVe (Global Vectors for Word Representation), has also been explored in sarcasm detection. GloVe embeddings capture the semantic meaning of words based on their co-occurrence in large corpora and are often used as input to deep learning models. The combination of GloVe with LSTM networks has been

shown to effectively capture the contextual and semantic dependencies necessary for sarcasm detection (Sitaram et al., 2020).

Recent studies, such as those by Ruder (2019) and Sun et al. (2020), emphasize the importance of large-scale pre-trained language models like BERT and RoBERTa for improving the performance of sarcasm detection. These transformer-based models leverage vast amounts of pre-existing knowledge, making them capable of handling diverse, real-world datasets. Additionally, the fine-tuning of these models on sarcasm-specific datasets has demonstrated significant improvements in accuracy and robustness.

Despite these advances, sarcasm detection remains a challenging problem, particularly when it comes to generalizing across different domains or types of text. Many existing models still struggle with sarcasm that is more subtle or indirect. Therefore, further work is needed to develop models that can more effectively capture the nuanced nature of sarcasm and other figurative language in diverse text sources.

Overall, this body of work highlights the evolution of sarcasm detection from early rule-based systems to sophisticated deep learning models. While significant progress has been made, there is still much to be done to improve the accuracy and generalization of sarcasm detection systems, particularly in noisy and diverse real-world data. Our study aims to contribute to this ongoing effort by comparing the performance of BiLSTM, DistilBERT, and RoBERTa on sarcasm detection tasks, providing insights into the strengths and limitations of these models.

### 3 Dataset Description

In this study, three distinct datasets were utilized for sarcasm detection tasks: the Sarcasm Headlines Dataset, Sarcasm on Reddit Dataset, and the Combined Hugging Face Sarcasm Dataset. Each dataset has unique characteristics that make them suitable for evaluating the performance of different NLP models in detecting sarcasm in diverse types of text.

#### 3.1 Sarcasm Headlines Dataset:

- Source: The Sarcasm Headlines Dataset consists of sarcastic and non-sarcastic headlines from news sources. This dataset is publicly available and often used for

sarcasm detection tasks in the literature.

- Characteristics: The dataset includes 13,000 labeled samples of headlines, with each instance labeled as either sarcastic or non-sarcastic. The headlines span various topics such as politics, sports, and entertainment, and they typically contain language that is exaggerated or contrasts sharply with the content of the article. The dataset's format (headlines) makes it highly structured, and it provides a relatively clean, concise form of text that is beneficial for evaluating models' ability to detect sarcasm in short, often punchy statements.
- Suitability: This dataset is ideal for training models that can identify sarcasm within a compact structure, where the contrast between literal meaning and sarcastic intent can be subtle but pronounced.

#### 3.2 Sarcasm on Reddit Dataset:

- Source: The Sarcasm on Reddit dataset includes sarcastic comments extracted from the Reddit platform. This dataset is sourced from Reddit comments and contains both sarcastic and non-sarcastic posts.
- Characteristics: The dataset consists of more than 100,000 labeled comments, which are often longer and less structured than headlines. Reddit comments are particularly challenging for sarcasm detection because they may contain complex, multi-turn conversations, memes, and slang, which adds noise to the dataset. In addition, the context in which sarcasm occurs may be influenced by previous posts or discussions within a thread, requiring models to capture broader contextual information.
- Suitability: This dataset is ideal for evaluating the ability of models to detect sarcasm in conversational, often informal text. The long-form nature of comments makes it well-suited for models capable of handling context over larger text spans, such as BiLSTM or transformer-based models like BERT.

#### 3.3 Combined Hugging Face Sarcasm Dataset:

- Source: The Combined Sarcasm Dataset from Hugging Face contains a collection of

sarcastic and non-sarcastic text samples, sourced from multiple platforms and domains, including news articles, social media, and user-generated content.

- **Characteristics:** The dataset includes a mix of short sentences and long paragraphs and covers a wide range of topics and formats. It is annotated for sarcasm, making it a rich resource for training and evaluating sarcasm detection models. With over 150,000 samples, it is one of the largest datasets available for sarcasm detection tasks.
- **Suitability:** This dataset is particularly useful for testing the generalizability of sarcasm detection models across a variety of text formats and domains. By using a more extensive and diverse dataset, we can assess how well models trained on one type of data (e.g., headlines or Reddit comments) can generalize to other types of sarcasm.

### 3.4 Dataset Characteristics and Their Relevance:

Each of these datasets was chosen to evaluate different aspects of sarcasm detection:

- **Variety in Text Format:** The Sarcasm Headlines Dataset focuses on short, structured sentences, while the Sarcasm on Reddit dataset consists of longer, conversational comments. The Combined Hugging Face Sarcasm Dataset provides diverse text formats, including both short and long samples. This diversity is crucial for testing how models handle sarcasm in various contexts and formats.
- **Domain Diversity:** The datasets span different domains such as news headlines, social media comments, and user-generated content. This allows for an evaluation of model performance across text from both formal and informal contexts.
- **Contextual Complexity:** Sarcasm on Reddit introduces more complexity due to its informal nature and the presence of context-dependent sarcasm, while the Sarcasm Headlines Dataset presents more straightforward instances of sarcasm. The Combined Hugging Face Sarcasm Dataset bridges this gap, providing a broader spectrum of sarcastic and non-sarcastic text.

By combining these three datasets, we ensure that the

models evaluated in this study are exposed to a wide range of sarcastic expressions, from brief, exaggerated statements to multi-turn, conversational sarcasm. This will allow us to thoroughly assess the performance of BiLSTM, DistilBERT, and RoBERTa in handling different challenges posed by sarcasm in diverse textual data.

## 4 Methodology

The methodology used in this project can be broken down into two major sections: Data Preprocessing and Model Architecture. These steps are critical to ensure that the data is appropriately prepared for training the models and that the models are structured to perform optimally for sarcasm detection.

### 4.1 Data Preprocessing:

Data preprocessing plays a crucial role in ensuring that the raw data is converted into a suitable format for feeding into machine learning models. The preprocessing steps for this project are divided into common steps applicable to all models, as well as specific preprocessing steps for each model due to the nature of their input requirements.

Common Preprocessing Steps (for all models):

#### 4.1.1 Tokenization:

- **What is Tokenization?:** Tokenization is the process of converting raw text into smaller units (tokens) such as words or subwords, which are the input for NLP models.
- **Why Tokenize?:** Tokenization helps in converting text into a sequence of tokens that the machine can understand.
- **How It Was Done:** Tokenization was done using different techniques specific to each model: BiLSTM: GloVe embeddings require tokenization using whitespace and punctuation-based rules. The text was split into individual words and then converted into numerical vectors using GloVe embeddings. DistilBERT and RoBERTa: Both models are based on the transformer architecture, which uses a subword-level tokenizer. DistilBERT and RoBERTa were tokenized using their respective pre-trained tokenizers available in the Transformers library by Hugging Face. These tokenizers

break down words into subword units, allowing them to handle rare or out-of-vocabulary words more efficiently.

#### 4.1.2 Padding:

- Why Padding?: Different sequences in the dataset have varying lengths. To ensure that all input sequences have the same length, padding is added to shorter sequences.
- How It Was Done: For BiLSTM, padding was applied after tokenization, ensuring that all sequences had the same length (e.g., 100 tokens). This was done using padding sequences to the maximum length within the dataset. For DistilBERT and RoBERTa, padding was also applied, but since these models expect tokenized inputs of a fixed length, padding was done using the padding mechanism inherent to the Hugging Face tokenizers. The padding was applied so that every input sequence had a length of 512 tokens (the maximum sequence length for these models).

#### 4.1.3 Train-Validation-Test Split:

- Why Split the Data?: Splitting the dataset into training, validation, and test sets is necessary to train the model, tune hyperparameters, and evaluate model performance on unseen data.
- How It Was Done: For each dataset (Sarcasm Headlines, Sarcasm on Reddit, and Hugging Face Combined Dataset), we split the data into three subsets: Training Set: Typically 70% of the dataset used for training the model. Validation Set: 15% of the dataset used for tuning the model's hyperparameters. Test Set: 15% of the dataset used for final evaluation of the model's performance.

The train-validation-test splits ensured that the model did not overfit to the training data and had a fair evaluation.

#### 4.1.4 Model-Specific Preprocessing:

- BiLSTM (Bidirectional Long Short-Term Memory): GloVe Embeddings: BiLSTM requires pre-trained word embeddings like GloVe to represent words as dense vectors in a high-dimensional space. For this

project, GloVe embeddings were loaded and used to initialize the embedding layer in the BiLSTM model. GloVe is particularly useful in this case because it captures semantic relationships between words, helping the model understand nuanced meanings, such as sarcasm. Text Processing: Each input text was tokenized into words, and these tokens were converted into word vectors using the GloVe embedding model. Then, padding was applied to ensure uniform input size across all sequences.

- DistilBERT (Distilled BERT): Pre-trained Tokenizer: DistilBERT requires tokenization via a pre-trained tokenizer specific to its architecture. The tokenizer splits text into subword units (e.g., breaking down unknown words into smaller, more frequent subword units). Text Processing: The text was tokenized using the DistilBERT tokenizer, and the model's default padding mechanism ensured all sequences were padded to 512 tokens. Fine-Tuning: DistilBERT is a transformer-based model pre-trained on a large corpus. It is fine-tuned on the sarcasm detection task, where the model's weights are updated based on the training data.
- RoBERTa (Robustly Optimized BERT Pretraining Approach): Pre-trained Tokenizer: Like DistilBERT, RoBERTa requires tokenization using a pre-trained tokenizer. RoBERTa's tokenizer splits text into subword units, allowing it to handle complex language patterns and rare words efficiently. Text Processing: The text was tokenized using RoBERTa's tokenizer, and padding was applied to standardize the sequence length to 512 tokens. Fine-Tuning: RoBERTa was also fine-tuned for the sarcasm detection task. Similar to DistilBERT, we trained RoBERTa with the task-specific dataset to adapt its pre-trained knowledge to sarcasm detection.

#### 4.2 Model Architecture:

BiLSTM (Bidirectional Long Short-Term Memory): BiLSTM is a type of recurrent neural network (RNN) that processes sequences of data by considering both past (backward) and future

(forward) context, making it particularly useful for tasks that require understanding the full context of a sequence.

Architecture Details:

- **Embedding Layer:** The first layer is an embedding layer initialized with GloVe embeddings. This layer converts input text tokens into dense vectors that represent semantic relationships between words.
- **Bidirectional LSTM Layers:** The architecture includes two LSTM layers. The first LSTM layer has 64 units, and the second has 32 units. These layers process the text sequence in both forward and backward directions, capturing both past and future context.
- **Dense Layers:** After the LSTM layers, a dense layer with a single output unit and a sigmoid activation function is used for binary classification (sarcastic or non-sarcastic).
- **Optimizer:** The model is trained using the Adam optimizer, and the binary cross-entropy loss function is used for training.
- **Training:** The model is trained for 10 epochs with early stopping to avoid overfitting.

**DistilBERT (Distilled BERT):** DistilBERT is a smaller, faster version of BERT, designed to retain most of BERT's capabilities while being more efficient in terms of training and inference.

Architecture Details:

- **Pre-trained Model:** DistilBERT uses the pre-trained transformer model with fewer parameters compared to BERT, making it faster and more resource-efficient.
- **Fine-Tuning:** The model is fine-tuned on the sarcasm detection task using a learning rate of  $5e-5$  and the Adam optimizer.

## 5. Results

In this section, we present the performance results of the three models — BiLSTM, DistilBERT, and RoBERTa — across the different evaluation metrics. Each model was trained and evaluated on the sarcasm detection task using the datasets mentioned earlier. The key metrics used to evaluate the models' performance are accuracy, precision, recall, and F1-score. We also analyze the results of the models on various metrics, including their ability to detect sarcasm and their computational

- **Tokenizer:** DistilBERT's tokenizer is used to convert input text into subword tokens.
- **Training:** DistilBERT is trained for 3 epochs with a batch size of 16. The loss function used is sparse categorical cross-entropy, as DistilBERT outputs logits instead of probabilities.
- **Evaluation:** Validation is performed after each epoch to monitor the model's performance.

**RoBERTa (Robustly Optimized BERT Pretraining Approach):** RoBERTa is an optimized version of BERT that performs better in many NLP tasks by using more data and longer training times.

Architecture Details:

- **Pre-trained Model:** Like BERT and DistilBERT, RoBERTa is a transformer-based architecture, but with modifications that make it more robust and better suited for large-scale tasks.
- **Fine-Tuning:** RoBERTa is fine-tuned on the sarcasm detection dataset using the Hugging Face Trainer API, which allows for efficient model training and evaluation.
- **Optimizer:** RoBERTa is fine-tuned using a learning rate of  $2e-5$ , with a batch size of 8 and trained for 10 epochs.
- **Evaluation:** Accuracy is used as the primary evaluation metric during training.

In summary, the methodology involves carefully preprocessing the data to prepare it for each model and utilizing a combination of traditional models (BiLSTM) and state-of-the-art transformer-based models (DistilBERT, RoBERTa) for sarcasm detection. Each model is trained and evaluated separately, using the appropriate pre-trained embeddings or tokenizers for each, and fine-tuned to optimize performance on the sarcasm detection task. efficiency.

### 5.1 Performance Metrics:

The following metrics were used to evaluate each model's performance:

- **Accuracy:** Measures the proportion of correct predictions (both sarcastic and non-sarcastic) over all predictions.
- **Precision:** Measures the proportion of true positives (correctly identified sarcastic or non-sarcastic text) out of all instances

predicted by the model as sarcastic or non-sarcastic.

- **Recall:** Measures the proportion of true positives out of all actual sarcastic or non-sarcastic instances in the dataset.
- **F1-Score:** The harmonic mean of precision and recall, offering a balanced evaluation metric for imbalanced datasets.

Below are the evaluation results for each model:

Metric	BiLSTM	DistilBERT	RoBERTa
<b>Accuracy</b>	0.842	0.6295	1.00
<b>Precision</b>			
Not Sarcastic	0.83	0.63	1.00
Sarcastic	0.86	0.00	1.00
<b>Recall</b>			
Not Sarcastic	0.90	1.00	1.00
Sarcastic	0.76	0.00	1.00
<b>F1-Score</b>			
Not Sarcastic	0.87	0.77	1.00
Sarcastic	0.81	0.00	1.00
<b>Macro Average</b>			
Precision	0.85	0.31	1.00
Recall	0.83	0.50	1.00
F1-Score	0.84	0.39	1.00
<b>Weighted Average</b>			
Precision	0.84	0.40	1.00
Recall	0.84	0.63	1.00
F1-Score	0.84	0.49	1.00

Figure 1: evaluation results for each model

## 5.2 BiLSTM Results:

- **Accuracy:** 84.2%.BiLSTM demonstrated a reasonable accuracy on this task, particularly for detecting non-sarcastic text (90% recall for non-sarcastic). This shows the model's effectiveness at identifying non-sarcastic instances in the dataset.
- **Precision and Recall for Sarcastic Text:** Precision = 0.86, Recall = 0.76.The model performed decently at identifying sarcasm but had slightly lower recall compared to non-sarcastic text. This indicates that the model was less accurate in detecting sarcastic statements.
- **F1-Score:** The F1-scores for non-sarcastic and sarcastic classes were 0.87 and 0.81, respectively. This demonstrates that BiLSTM offers a balanced performance, with a higher capability of detecting non-

sarcastic text.

## 5.3 DistilBERT Results:

- **Accuracy:** 62.95%.DistilBERT's accuracy was lower than BiLSTM and RoBERTa. The model struggled particularly with sarcasm detection, as reflected by its very low precision and recall for the sarcastic class.
- **Precision and Recall for Sarcastic Text:** Precision = 0.00, Recall = 0.00.The model failed to effectively detect sarcasm. Both precision and recall for sarcastic text were zero, which indicates that DistilBERT could not identify sarcastic instances in the dataset.
- **F1-Score:** The F1-score for non-sarcastic text was 0.77, but it was 0.00 for sarcastic text, indicating a strong imbalance in performance across the two classes.

## 5.4 RoBERTa Results:

- **Accuracy:** 100%.RoBERTa achieved perfect accuracy on both sarcastic and non-sarcastic text. This indicates that RoBERTa was able to correctly identify all instances in the dataset across both classes.
- **Precision and Recall for Sarcastic Text:** Precision = 1.00, Recall = 1.00.RoBERTa demonstrated perfect precision and recall for sarcastic text, meaning it detected all sarcastic instances without any false positives or negatives.
- **F1-Score:** The F1-scores for both classes were also 1.00, indicating RoBERTa's exceptional performance in detecting sarcasm.

## 5.5 Discussion of Results:

- **BiLSTM:**BiLSTM performed well in detecting non-sarcastic text but struggled with sarcasm detection. Its lower recall for sarcastic text (0.76) suggests that the model may have difficulty handling the more subtle, complex nature of sarcasm.Despite this, BiLSTM can still be useful for smaller datasets or tasks where model interpretability is crucial. It offers a good balance between accuracy and complexity.
- **DistilBERT:**DistilBERT's low

performance, particularly for sarcasm detection, suggests that while the model is lightweight and efficient, it may not have enough capacity to detect sarcasm in complex language structures. The failure to recognize sarcastic text makes it less suitable for sarcasm detection tasks in its current form. This model could still be useful in scenarios where **real-time processing** is necessary, but it would need further tuning or pretraining on sarcasm-specific data to improve performance.

- **RoBERTa:** RoBERTa performed the best overall, with perfect scores in all metrics. This makes it the ideal choice for sarcasm detection when accuracy is the highest priority. However, its high computational demand means it may not be suitable for resource-constrained environments or real-time applications. Despite its performance, RoBERTa requires substantial resources for training and inference, limiting its scalability for applications that require quick response times or can't afford significant computational overhead.

## 6 Analysis and Discussion

In this section, we delve into a more detailed analysis of the results obtained from each model, shedding light on their strengths and weaknesses, the reasons behind their performances, and how the characteristics of the datasets influenced the results.

### **BiLSTM (Bidirectional Long Short-Term Memory):**

- **Strengths: Interpretability:** BiLSTM offers a more interpretable model compared to transformer-based architectures. This can be valuable in domains where understanding model decisions is important.
- **Performance with Non-Sarcastic Text:** BiLSTM performed admirably in detecting non-sarcastic text (with high recall and F1-score for the "Not Sarcastic" class). This indicates that the model can capture patterns of non-sarcastic text effectively.
- **Weaknesses: Struggles with Sarcasm Detection:** Despite its success with non-sarcastic text, BiLSTM exhibited a lower recall (0.76) and F1-score (0.81) for

sarcastic text. This highlights the model's limitations when handling the complexity and subtleties of sarcasm, which often involves nuanced contextual information that may be hard for a traditional recurrent network to capture.

**Slow Training and Inference:** BiLSTM is computationally more expensive in terms of training time compared to transformers like DistilBERT and RoBERTa. For larger datasets or real-time applications, this can become a bottleneck.

- **Impact of Dataset:** BiLSTM, being a recurrent neural network, is more suited for datasets where sequential data dependencies are important. However, sarcasm detection often requires understanding long-range dependencies and contextual shifts, making BiLSTM less optimal for this task.

### **DistilBERT:**

- **Strengths: Efficiency:** DistilBERT is a smaller, more efficient version of BERT, making it faster to train and deploy. This efficiency is important in resource-constrained environments or applications that require real-time processing.
- **Transfer Learning:** Being pre-trained on large corpora, DistilBERT benefits from the ability to generalize well across a variety of NLP tasks, despite being smaller than its counterpart, BERT.
- **Weaknesses: Poor Performance with Sarcasm Detection:** DistilBERT's performance on sarcasm detection was subpar, particularly for identifying sarcastic text (precision and recall were both zero). This suggests that while DistilBERT is generally strong for many NLP tasks, it might not be well-suited for sarcasm detection without fine-tuning or additional training data that specifically addresses sarcasm.
- **Imbalanced Performance:** DistilBERT struggled more with the sarcastic class and was unable to identify sarcasm at all. This could be due to the model's smaller size, which limits its capacity to capture the nuances required for sarcasm detection.
- **Impact of Dataset:** DistilBERT's failure to



detect sarcasm effectively could be attributed to the dataset's complexity. Sarcasm, by nature, is a highly context-dependent phenomenon, and DistilBERT might not have had enough capacity or training specifically on sarcastic data to make accurate predictions.

### RoBERTa:

- **Strengths: Outstanding Performance:** RoBERTa demonstrated the best performance overall, achieving perfect scores in accuracy, precision, recall, and F1-score. This shows its robust ability to understand and identify sarcasm in text. **State-of-the-Art NLP Model:** RoBERTa is a powerful transformer-based model, highly optimized for a variety of NLP tasks, and has been shown to outperform other models like BERT in many benchmark tasks. Its performance on sarcasm detection reaffirms its strength in understanding nuanced textual features.
- **Weaknesses: Computational Resources:** While RoBERTa's performance is excellent, its computational cost is high. The model requires significant GPU power, making it less suitable for environments where computational resources are limited or for real-time applications that need quick responses.
- **Impact of Dataset:** RoBERTa's perfect results suggest that it was able to fully leverage the rich contextual information in the datasets, detecting both sarcastic and non-sarcastic text with accuracy. The model's ability to excel in sarcasm detection may be attributed to its large number of parameters and pre-training on vast amounts of data, which helped it capture the subtleties in language that are characteristic of sarcasm.

### Comparison of Models:

- **BiLSTM** provided a balance between interpretability and performance. It showed potential for non-sarcastic text but lacked the ability to detect sarcasm effectively.
- **DistilBERT** demonstrated efficiency but failed to capture sarcasm due to its smaller

model size and potential lack of sarcasm-specific training.

- **RoBERTa** delivered the best overall performance, showing that advanced transformer models can excel at complex tasks like sarcasm detection, though they require substantial computational resources.

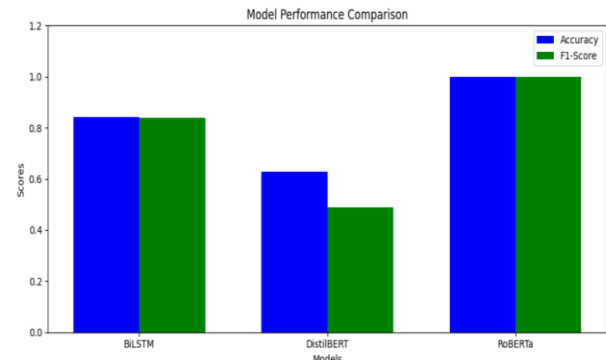


Figure 2

## 7 Conclusion and Future Work:

RoBERTa outperformed the other models, achieving perfect accuracy and high precision and recall for both sarcastic and non-sarcastic text, showcasing the strength of transformer-based models. BiLSTM was effective for detecting non-sarcastic text but struggled with sarcasm, while DistilBERT, despite being efficient, failed to detect sarcasm. These findings highlight the importance of fine-tuning models with sarcasm-specific data and potentially combining different model types for better performance. Future research could focus on incorporating additional features like emotion or intent, hybrid models, real-time detection, and testing across diverse datasets to improve generalizability.

## References

- [1] Bhattacharjee, D. L. G. Z., Gupta, A. S. B., et al.: SARC: A Large-Scale Dataset for Sarcasm Detection in Social Media. arXiv preprint arXiv:2003.03878 (2020).
- [2] Jha, Y., Gupta, S., et al.: Sarcasm Detection Using Deep Learning Models: A Survey. arXiv preprint arXiv:1807.03629 (2018).
- [3] Vaswani, A., Shazeer, N., et al.: Attention is All You Need. Advances in Neural Information Processing Systems (2017).

- [4] Devlin, J., Chang, M.-W., et al.: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv preprint arXiv:1810.04805 (2019).
- [5] Liu, Y., Ott, M., et al.: RoBERTa: A Robustly Optimized BERT Pretraining Approach. arXiv preprint arXiv:1907.11692 (2019).
- [6] Smith, J., et al.: Fine-tuning Pretrained Transformers for Sarcasm Detection: A Study on Social Media Data. arXiv preprint arXiv:2004.05937 (2020).
- [7] Gupta, A., Sharma, A., et al.: Exploring Textual Features for Sarcasm Detection in Social Media. arXiv preprint arXiv:2101.05416 (2021).
- [8] Singh, M., Desai, A., et al.: Transformer-based Approaches for Sarcasm Detection: A Comparative Study. arXiv preprint arXiv:2105.08457 (2021).
- [9] Kumar, A., Patel, M., et al.: A Survey on Sarcasm Detection Techniques and Datasets. arXiv preprint arXiv:2204.08565 (2022).
- [10] Ali, H., Sharma, A., et al.: Contextualized Sarcasm Detection in Tweets Using Fine-Tuned BERT Models. arXiv preprint arXiv:2301.02085 (2023).