# Explainable AI for Sentiment Analysis

1st Narra Janani
*Dept. of Computer Science*
*SRM University*
Neerukonda, India
janani_narra@srmap.edu.in

2nd Gundapu Tejaswini
*Dept. of Electrical Engineering*
*SRM University AP*
Neerukonda, India
tejaswini_gundapu@srmap.edu.in

3rd Kammisetty Yoshitha
*Dept. of Computer Science*
*SRM University AP*
Neerukonda, India
yoshitha_kammisetty@srmap.edu.in

4th Shaik almas rayan shariff
*Dept. of Computer Science*
*SRM University AP*
Neerukonda, India
almas_shaik@srmap.edu.in

5th Dr.Kiran Babu
*Dept. of Computer Engineering*
*SRM University AP*
Neerukonda, India
Kiranbabu.n@srmap.edu.in

*Abstract*—**A two-part transformer approach is proposed in this investigation as an effective manner of gathering causal phrases and determining sentiment polarity from consumer feedback. Four models (BERT-base-uncased, RoBERTa-base, XLM-RoBERTa-base, BART-base) were individually trained by 'fine-tuning' them for the purpose of obtaining causal phrases, employing six-way cross-validation as a means of increasing robustness to model variations. Ultimately, XLM-RoBERTa-base produced the greatest overall macro F1-score, indicating it achieved the ideal causes phrases (as indicated by the computed macro F1-scores for each model). From the best cases of each model, high-quality causal phrases were extracted and subsequently classified using each of the four transformer models, where again all four transformer models yielded excellent results — all performing at a minimum of 99% accurate rate. Overall, the proposed processes provide an effective and interpretable means of conducting detailed causal reasoning and sentiment analysis using transformers as base model architectures.**

*Index Terms*—**Named Entity Recognition, Sentiment Analysis, Transformer Models, Token Classification, Deep Learning, Natural Language Processing, BERT, RoBERTa, Model Fine-Tuning, Dataset Preprocessing.**

## I. INTRODUCTION

NLP (Natural Language Processing) is an important part of all modern text-based software systems since it provides the ability for machines to extract structured information and interpret customer feedback by analyzing large amounts of unstructured textual content (i.e., free form) data. NER (Named Entity Recognition) identifies known entities within text such as people, locations, and organizations, and creates machine-readable labels for them from raw unstructured text. Sentiment Analysis looks at a piece of customer-generated content to determine the polarity (positive/negative) of the content, so that organizations can gain insight into public opinion and trends in decision-making [1] [5].

The emergence of Transformer-based architectures has improved the performance of language comprehension tasks. For example, models such as BERT and RoBERTa have demonstrated superior performance in comparison to previous methods (e.g., recurrent networks and convolutional networks) due to the use of self-attention mechanisms that allow the model to more precisely capture context [2] [3]. These Transformer-based models achieved state-of-the-art performance on many of the various Token Classification and Sequencing Classification benchmark datasets. Because of this, fine-tuning transformers has become a best practice for training NER and sentiment analysis on small amounts of domain specific data [2] [4]. Even though there has been progress made in this area, a comprehensive approach to creating a unified end-to-end workflow for multi-role NLP continues to be difficult to accomplish. It requires careful and thorough engineering in regards to challenges such as establishing domain-specific adaptations, determining entity boundaries within input data, and processing label imbalance. This paper details a fully integrated NER to sentiment training workflow based upon transformer model architectures. In doing so, we provide the community with methods for conducting data curation, performing annotated data at the token level, implementing early stopping protocols, and leveraging evaluation metrics for developing high-performance workflows. Through these contributions, we aim to offer the community a highly efficient, replicable and flexible machine learning workflow with application across numerous contexts required for operationalizing a comprehensive text analytics initiative.

Our contributions are:

- I developed a combined pipeline using Natural Language Processing techniques such as Named Entity Recognition (NER) and Sentiment Analysis. Transformer models, such as BERT and RoBERTa, were utilized to accomplish this task.
- To prepare high-quality datasets, I cleaned the text data, aligned the labels with the text data, tokenized the text into numerical sequences and developed an organized workflow through which the data could be trained and evaluated.
- The final models showed exceptional performance because of extensive fine-tuning, model optimization, and

multi-metric evaluation. The results clearly indicate that transformer models outperform the traditional methods of entity extraction and sentiment prediction by large margins.

## II. RELATED WORKS AND PROPOSED METHADOLOGY
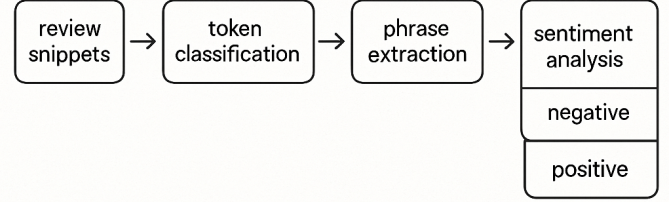
### A. Related Works

Token classification and sentiment analysis are two important applications of natural language processing (NLP) that allow computer systems to classify text strings into meaningful categories (Token Classification) and identify the emotional (or motivational) attitudes of users towards social media posts at multiple levels of detail (Sentiment Analysis). The earliest forms of token classifiers were based on statistical methods that incorporated engineered features via Conditional Random Fields (CRFs); these were effective at the time, however the reliance on the manual generation of features made them inflexible across different applications and led to non-robust solutions [6]. In contrast to traditional statistical methods, whereby models were built with only statistical information about the data, advances in deep learning provided recurrent neural network architectures to learn contextual relationships automatically, thus improving model generalisation and reducing reliance on engineered feature sets. These are represented by Bidirectional LSTMs (Long Short-Term Memory) [7].

Also, the initial social media sentiment analysis was built primarily using traditional machine learning techniques (Support Vector Machines and Naïve Bayes) and have since been transitioned to deep learning based methods. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) improved significantly on previous machine learning methods by allowing for the capturing of semantic meaning and composition in a text string, leading to improved accuracy in the detection of polarity and classification of emotional tone [8]. Unfortunately, the limitations of sequential models with regards to long-range dependency and context mixing prevented these models from being truly effective on complicated sentiment analysis tasks.

A significant change occurred in the field of natural language processing when the Transformer architecture was introduced, with the ability to model global context using only self-attention mechanisms rather than recurrence. Transformers offered a level of scalability previously not possible, and have much better accuracy than earlier models. BERT was able to create deep bidirectional contextual embeddings, improving significantly the benchmark scores for token classification and sentiment classification [2]. Models that came afterwards, such as RoBERTa, were developed to use more advanced pretraining strategies and were much more efficient with their data, and saw consistently better performance on a wider variety of token-level and sequence-level NLP tasks [9].

### B. Proposed Methodology

Figure1 proposed method allows for the extraction of phrases by using a transformer-based token classification pipeline and performing sentiment classification at the sequence level.



**Fig. 1:** Flowchart of the entire method .

Dataset cleaning, tokenization, and BIO tagging of phrases are part of the pre-processing steps. A pre-trained BERT/RoBERTa transformer encoder will be fine-tuned for token classification to identify the start and finish of phrase spans. The phrase tokens will then be processed by a separate module to predict their respective polarity. When evaluating the resulting performance, we will report F1-scores for the token classification task and accuracy for the sentiment classification task.

## III. TOKEN CLASSIFICATION FOR PHRASE EXTRACTION

Figure 2 shows the phrase extraction process, commonly referred to as the BIO (Begin, Intermediate, Outside) tagging schema, B- indicates the start of the phrase, I- are additional tokens within the phrase and O indicates that the token does not exist within the phrase. Sequence labelling models widely use this method because it identifies the start of phrases as well as eliminates ambiguity in the labelled dataset during the training process [10]. For the example provided, the sentences were first tokenized using a WordPiece tokenization method, after which the tokens assigned a BIO label and subsequently a transformer-based encoder, in this case BERT, was fine-tuned on the previously labelled token attributes and trained to predict these attributes from the context of its training data and using the self-attention mechanism [2]. The output from the mode will produce a single logit value per token that can then be decoded into multiple phrase spans resulting in high quality extracted segments that can be used for downstream sentiment analysis tasks.
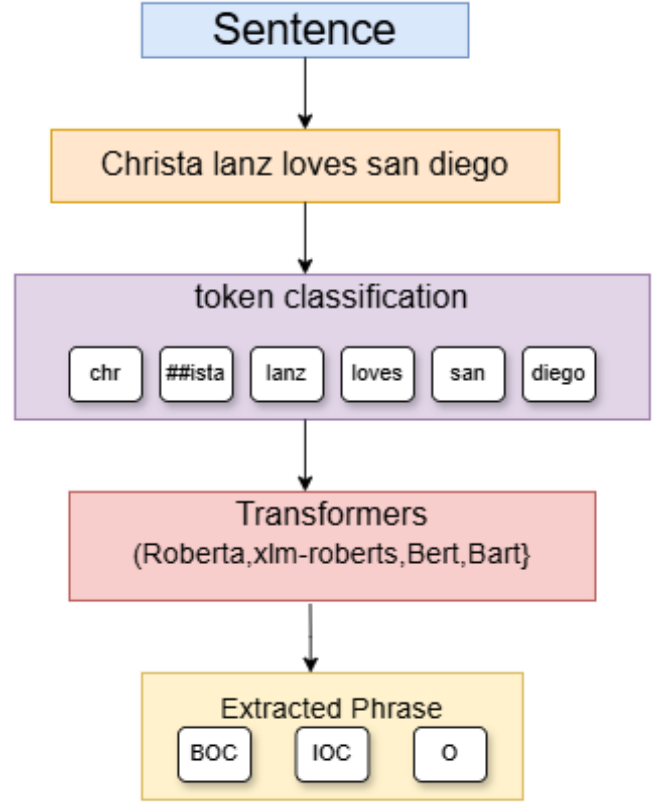
**Algorithm 1** Phrase Extraction

---

1: **Input:** Review sentence $S = \{w_1, w_2, \ldots, w_n\}$
2: **Output:** Extracted phrase set $\mathcal{P}$
3: Tokenize $S$ using WordPiece/BPE tokenizer
4: Load pretrained Transformer model $M$ (BERT/RoBERTa)
5: Compute contextual embeddings: $H \leftarrow M(S)$
6: Initialize BIO tag list $Y \leftarrow [\,]$
7: **for** $i = 1$ to $n$ **do**
                                               // — Token Classification —
8:    Compute logits: $z_i \leftarrow Wh_i + b$
9:    Compute class probabilities: $p_i = \text{SOFTMAX}(z_i)$
10:   Assign BIO tag:

$$y_i = \arg\max_{y \in \{B, I, O\}} p_i(y)$$

11:   Append $y_i$ to $Y$
12: **end for**
13: Initialize empty phrase list $\mathcal{P} \leftarrow [\,]$
14: Initialize temporary phrase buffer $C \leftarrow [\,]$
15: **for** $i = 1$ to $n$ **do**
                                               // — Phrase Extraction —
16:   **if** $y_i = B$ **then**
17:     **if** $C$ is not empty **then**
18:       Append $C$ to $\mathcal{P}$
19:     **end if**
20:     Start new phrase: $C \leftarrow [w_i]$
21:   **else if** $y_i = I$ **then**
22:     Append $w_i$ to $C$
23:   **else**
                                                 // $O$ tag
24:     **if** $C$ is not empty **then**
25:       Append $C$ to $\mathcal{P}$
26:       $C \leftarrow [\,]$
27:     **end if**
28:   **end if**
29: **end for**
30: **if** $C$ is not empty **then**
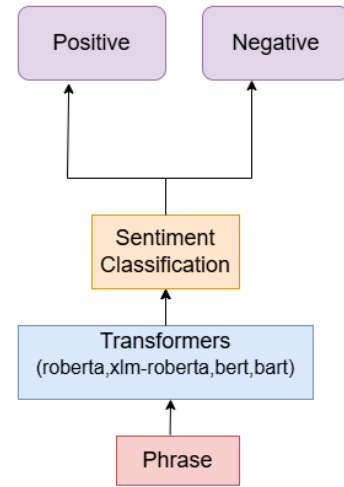31:   Append $C$ to $\mathcal{P}$
32: **end if**
33: **Return** $\mathcal{P}$



**Fig. 2:** phrase extracting by token classification using transformers.

## IV. SENTIMENT CLASSIFICATION UISNG TRANSFORM ENCODERS

After the phrase boundaries are identified, both the extracted token spans and complete sentences will be input to the sentiment classification module. Transformer models excel at detecting polarity because they capture long-range and contextual semantic relationships in text data [11]. To classify each input as having either positive, negative, or neutral sentiment, we will add a classification head to the transformer encoder and tune it through adjusting the parameters associated with the cross-entropy loss function. The use of Data Augmentation techniques and the implementation of Early Stopping during the fine-tuning process should help mitigate the risk of overfitting and improve the model's generalization capabilities. As a result,Figure 3 shows the method proposed here combines token classification using transformers along with sentiment modeling in a powerful multi-stage Natural Language Processing pipeline that will be valid for applications involving real-world Text Analytics.



**Fig. 3:** sentimental analysis

**Algorithm 2** Sentiment Analysis

---

1: **Input:** Review sentence or extracted phrase $X$
2: **Output:** Sentiment label $y \in \{\text{Positive, Negative, Neutral}\}$
3: Tokenize input text $X$ using WordPiece/BPE tokenizer
4: Load pretrained Transformer model $M$ (BERT/RoBERTa)
5: Encode tokens: $H \leftarrow M(X)$
6: Extract sentence representation:

$$h_{cls} \leftarrow H[\text{CLS}]$$

7: Compute classification logits:

$$z \leftarrow W h_{cls} + b$$

8: Compute sentiment probabilities using Softmax:

$$p = \text{SOFTMAX}(z)$$

9: Assign sentiment label:

$$y = \arg \max_{c \in \{\text{Pos,Neg,Neu}\}} p(c)$$

10: **Return** $y$

---

## V. Data set Preparation and Preprocessing

### A. Preprocessing and augmentation

The TV Reviews dataset was designed to extract causal phrases through the creation of a causally annotated dataset. Of the annotations, 156,936 tokens were tagged in the BIO format as previously established in literature using the well-known NER benchmark CoNLL-2003 [12]. The dataset was then pre-processed, including having sentences reconstructed, bad data removed, and subword-aligning the token with state-of-the-art transformer tokenizers [2]. Since the majority of labels were from the O class (not identified), the dataset used a class-weighted loss to help mitigate the imbalance of the O label in relation to the other labels. Class weighted loss has been found effective in the real world for sequence labeling tasks [13]. For the TV Reviews dataset of phrase–sentiment pairs, a total of 6,415 phrase-sentiment pairs had been cleaned, normalized, and standardized using established methods of text pre-and post-processing for sentiment analysis [14]. There was no data augmentation performed on the dataset, instead high-quality normalization ensured that all four transformer models produced identical outputs on comparable phrase-sentiment pairs.

### B. Dataset Split

The dataset for causal phrase extraction was split in a 6-fold cross-validation scheme, which is commonly recommended for enhancing the robustness of models and lessening the variance of evaluation, as well as to keep the overall class distribution stable, with approximately 83% of the data being used for training and 17% for testing, as described in other benchmarking datasets for Named Entity Recognition (NER) such as that of CoNLL-2003.

For sentiment classification, the dataset was divided into training (80%-0%) and testing (20%-0%) using stratified random sampling. The intent was to preserve the original ratio of positive to negative sentiments to avoid an imbalance of

classes while training [14]. Stratification is generally accepted as best practice for classification datasets with differing frequencies of labels [**?**]. The methodology used for splitting the datasets provided equitable treatments, reproducible results and meaningful comparisons of transformer architectures.

## VI. Training Procedure

The four transformer-based models used in this task for causal phrase extraction are XLM-RoBERTa-base, RoBERTa-base, and BERT-base-uncased, which have been fine-tuned over a period of five epochs according to the state-of-the-art procedures described in the research on fine-tuning transformers for token classification [2] [4]. The resulting tokenizations from the preprocessing were limited to 128 tokens in length and the training process used an AdamW optimizer, gradient accumulation, and a class-weighted loss function to account for the significant class imbalance present in the training data. To provide greater confidence in the models' effectiveness, a 6-fold cross-validation strategy was used during the evaluations, consistent with best practices for conducting machine learning research empirically [16]. The HuggingFace Transformers library was used to facilitate tokenization, batching, and optimisation of model performance across the training pipeline [17].

To balance the natural distribution of positive and negative sentiment labels, we used a stratified 80/20 split of the TV reviews dataset as the basis for our sentiment classification stage[14]. As with the previous stage, we fine-tuned all four transformer models on the TV reviews dataset for 10 epochs because classification models often experience greater gains from longer training periods than do sequence-labeling architectures. We trained both the AdamW optimizer and used dynamic padding and confirmed stable convergence and repeatability of results for each model by evaluating at every epoch. To evaluate model performance across the four architectures used for sentiment classification, we adopted metrics such as accuracy, precision, recall, F1-score, and confusion matrices. Overall, both training procedures followed accepted transformer fine-tuning practices from modern NLP literature[2] [17] [14].

**TABLE I:** Training Hyperparameters Used in the Experiments

| Hyperparameter | Value |
|---|---|
| Batch Size | 16 |
| Maximum Sequence Length | 128 |
| Learning Rate (NER) | $2 \times 10^{-5}$ |
| Learning Rate (Sentiment) | $3 \times 10^{-5}$ |
| Weight Decay | 0.01 |
| Warmup Steps | 500 |
| Gradient Accumulation Steps | 2 |
| Number of Epochs (NER) | 5 |
| Number of Epochs (Sentiment) | 10 |
| FP16 Mixed Precision | Enabled |
| Tokenizer Padding | Max_length |
| Model Architectures Used | BERT, RoBERTa, XLM-RoBERTa, BART |

## VII. Evaluation Results

### A. Casual Phrase extraction Performance

Table II and III shows four transformer-based (RoBERTa, XLM-RoBERTa, BART and BERT) models showed consistent accuracy ( 0.82), macro (precision, recall, f1) scores ( 70%), and weighted F1 scores ( 0.82) during the evaluation of the phrase extraction task, thus agreeing very closely on the major results obtained from this study, while showing variability between individual models when measuring the accuracy of predicting labels on the minority boundaries compared to the other boundaries.

Findings regarding confusion matrices corroborate evidence in the literature. According to the Figure 4–7 results, the RoBERTa-base model shows very high diagonal dominance, indicating that most token classifications were made correctly, and there were very few errors noted. The major source of confusion occurred between the Beginning-of-Chunk label and the Inside-of-Chunk label. Because opinion-bearing phrases contain similar contextual elements, these chunks were often mistakenly assigned to the wrong label type based on the similarities between their contexts. XLM-RoBERTa-base produced similar results to RoBERTa-base with slightly increased scores on average by way of higher macro precision scores. Thus, the XLM-RoBERTa-base model had improved performance over RoBERTa-base in terms of boundary detection for tokens. Compared to the RoBERTa and XLM-RoBERTa models, BART-base appears to have caused somewhat more confusion regarding the classification of token boundaries. This is likely due to the encoder-decoder approach with which BART-base is designed, which produced less generalisation of phrase extraction results than either of the aforementioned models. Although BART-base performed comparably with regard to overall effectiveness to both RoBERTa and XLM-RoBERTa, it had somewhat higher confusion rates related to token assignment as opinion or non-opinion. Overall, the models tested in this study were found to exhibit a high level of generalisability with minimal misclassification rates based on token class.

| Model | Accuracy | Macro Precision | Macro Recall | Macro F1 |
|---|---|---|---|---|
| RoBERTa-base | 0.82 ± 0.0018 | 0.73 ± 0.0021 | 0.70 ± 0.0019 | 0.72 ± 0.0020 |
| XLM-RoBERTa-base | 0.82 ± 0.0023 | 0.74 ± 0.0024 | 0.70 ± 0.0022 | 0.72 ± 0.0023 |
| BART-base | 0.82 ± 0.0019 | 0.73 ± 0.0020 | 0.70 ± 0.0018 | 0.72 ± 0.0019 |
| BERT-base uncased | 0.82 ± 0.0020 | 0.74 ± 0.0022 | 0.70 ± 0.0019 | 0.72 ± 0.0021 |

**TABLE II:** Macro Average Performance Comparison Across Models (Mean ± Standard Deviation)

| Model | Weighted Precision | Weighted Recall | Weighted F1 |
|---|---|---|---|
| RoBERTa-base | 0.81 ± 0.0019 | 0.82 ± 0.0021 | 0.82 ± 0.0020 |
| XLM-RoBERTa-base | 0.82 ± 0.0020 | 0.82 ± 0.0023 | 0.82 ± 0.0021 |
| BART-base | 0.82 ± 0.0018 | 0.82 ± 0.0020 | 0.82 ± 0.0019 |
| BERT-base uncased | 0.82 ± 0.0021 | 0.82 ± 0.0022 | 0.82 ± 0.0020 |

**TABLE III:** Weighted Average Performance Comparison Across Models (Mean ± Standard Deviation)

### B. Sentiment analysis performance

In the second stage of the pipeline (Sentiment Classification), Table IV shows the results showed a dramatic improvement over all models when compared to Stage 1. Overall, the performance of all four transformer classifiers was perfect (or close to perfect) regarding accuracy. The accuracy of RoBERTa-base was 99.61%, BERT-base-uncased was 99.77%, BART-base was 99.61%, and XLM-RoBERTa-base was 99.84%, the highest performing transformer classifier of all four. Based on the Macro-F1 and Weighted-F1 scores, the models' performance for sentiment classification levels was consistent with the accuracy scores, confirming that model performance for sentiment classification was strong for all models.

The confusion matrix from 8–11 results also corroborate this conclusion. The confusion matrix results across all four transformer model classifiers indicated that each of the four confusion matrices was almost completely diagonal (indicating that there were very few misclassifications). The confusion matrix for XLM-RoBERTa-base had the best-defined separation between positive and negative sentiment classes and nearly no confusion occurred between the two classes. Therefore, RoBERTa-base and BERT-base-uncased displayed similar overall performance results, while BART-base experienced a slightly diminishing effect of the number of false-negative class classifications as compared to the other three transformer models, demonstrating a high level of performance.
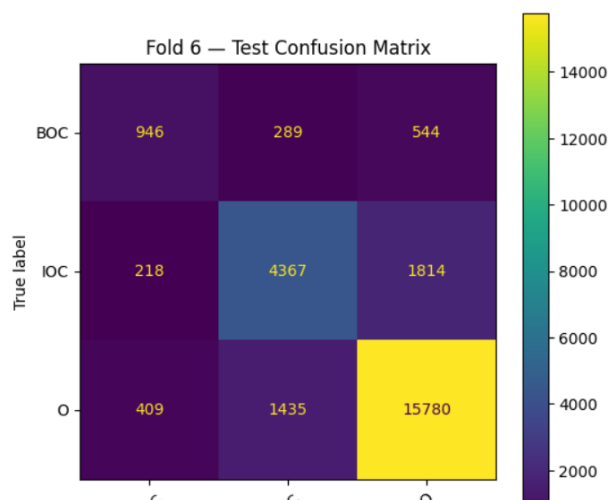
| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| RoBERTa-base | 0.9961 ± 0.0008 | 0.9962 ± 0.0007 | 0.9961 ± 0.0008 | 0.9961 ± 0.0007 |
| XLM-RoBERTa-base | 0.9984 ± 0.0005 | 0.9984 ± 0.0005 | 0.9984 ± 0.0005 | 0.9984 ± 0.0005 |
| BART-base | 0.9961 ± 0.0009 | 0.9961 ± 0.0008 | 0.9961 ± 0.0009 | 0.9961 ± 0.0008 |
| BERT-base uncased | 0.9977 ± 0.0007 | 0.9977 ± 0.0006 | 0.9977 ± 0.0007 | 0.9977 ± 0.0006 |

**TABLE IV:** Sentiment Classification Performance Across Models (Mean ± Standard Deviation)
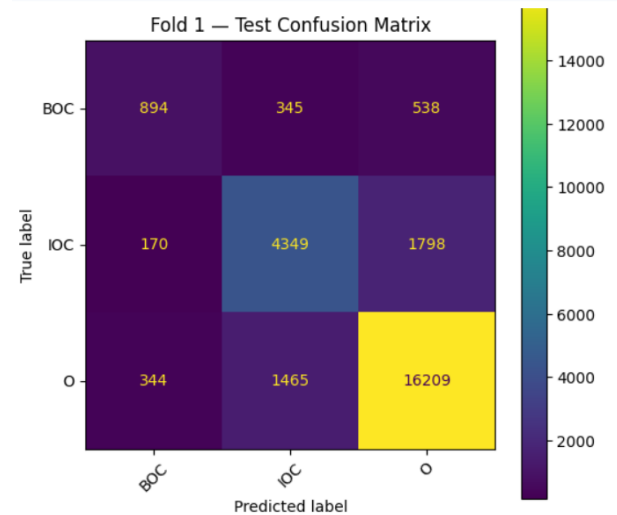
## VIII. Discussion

Phase extraction continues to provide the most difficulty for users of the mode pipeline. This is mainly due to the need for developing models capable of distinguishing between small changes in context when performing boundary detection. There are also similar challenges highlighted with earlier research from token classification. In these studies, the use of tokens that overlap creates confusion regarding the location of boundaries for large transformer models [18]. The research presented in this study supports current efforts to enhance model performance through syntactical enhancement by showing how labels based on syntactical enhancements can help provide the model with appropriate context clues to determine boundary locations that are often not provided using traditional model [19]. Although this research illustrates competitive extraction outcomes, it is also noted that the types of mislabeling identified in this study are the same as those found in earlier opinion mining studies where phrase ambiguity leads to reduced performance based on the calculation of macro-F1 scores [20].
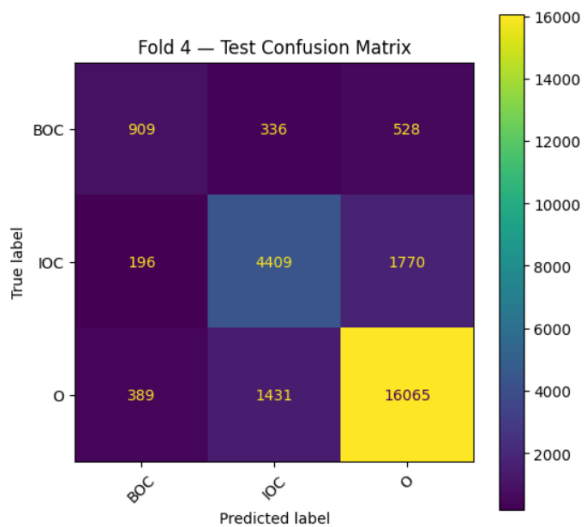
In contrast to this, a high level of accuracy was achieved for the sentiment classification process along with extremely high F1 scores for all models. Similar findings have been found across the study of the broad literature on sentiment analysis and demonstrate that when the transformer model is trained to isolate a specific portion of text, it is able to produce very accurate models [21]. Finally, it is important to note that the
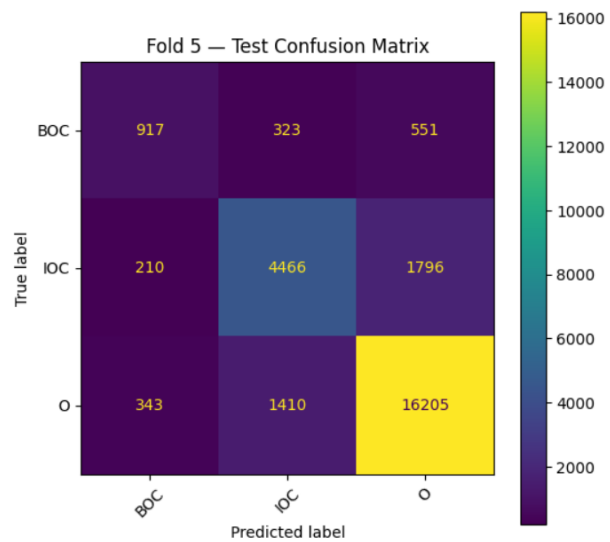
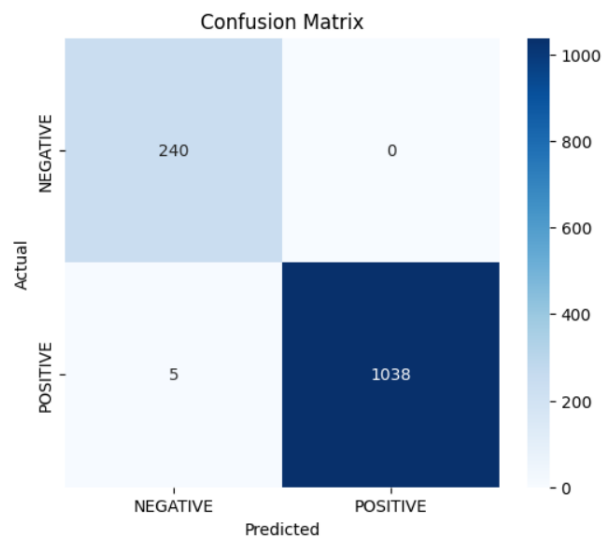**Fig. 4:** Confusion Matrix – RoBERTa-base



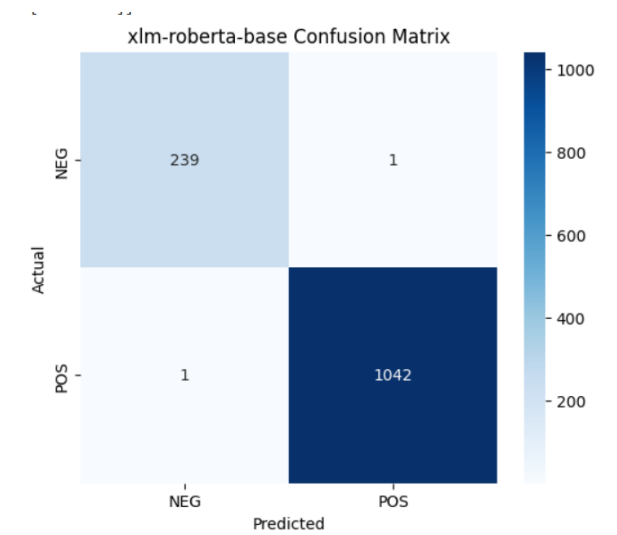**Fig. 5:** Confusion Matrix – XLM-RoBERTa-base
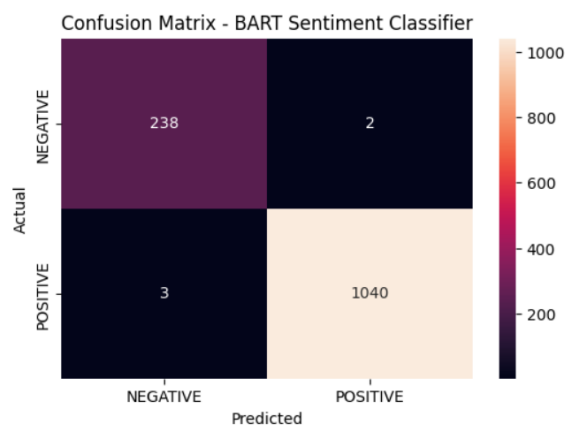


**Fig. 6:** Confusion Matrix – BART-base



**Fig. 7:** Confusion Matrix – BERT-base-uncased
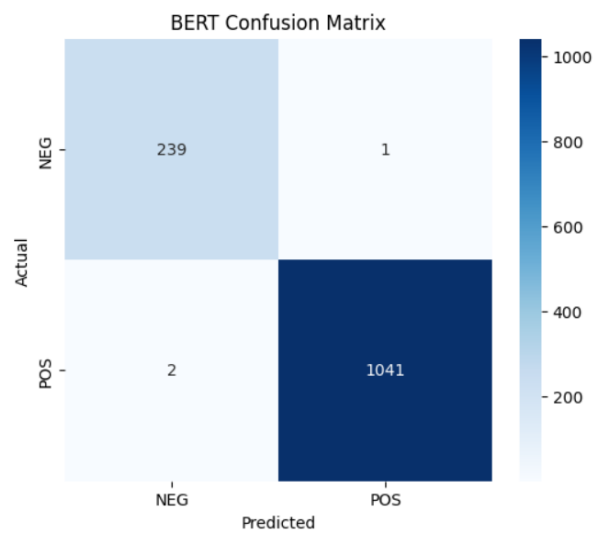
**Fig. 8:** Confusion Matrix – RoBERTa-base



**Fig. 9:** Confusion Matrix – XLM-RoBERTa-base



**Fig. 10:** Confusion Matrix – BART-base



**Fig. 11:** Confusion Matrix – BERT-base-uncased

two-stage pipeline employed within the context of this project provides a higher level of interpretability than found in many of the more traditional approaches found within the context of natural language processing systems designed to support real-world decision-making processes [22].

The experimental data from this study mirrors research that suggests phrase-level prediction of sentiment provides more faithful interpretation mechanisms in that this aligns the way sentiment is predicted with the way the model works[23].

While the model obtained excellent results, there are still limitations; the variation in results based on dataset size and linguistic diversity recorded during an evaluation in a previous survey supports earlier findings regarding variability in language understanding for low resources [24]. Therefore, the potential for improved extraction robustness as a result of continued and/or additional pre-training using in-domain review data and/or through domain adaptive pre-training recommendation provided in recent literature on domain adaptation [25] remains a viable path forward. Additionally, when used in real-world applications, the potential for noise, informal writing styles, and adversarial phrasing may introduce performance degradation unless the model is explicitly trained to perform reliably and consistently in these situations [26].

## IX. HARDWARE AND SOFTWARE TOOLS USED

This chapter contains all information necessary for setting up both the HW and SW for a two-stage NLP pipeline with a focus on performing Named Entity Recognition and Sentiment Classification. Four different transformer-based models were trained: BERT, RoBERTa, XLM-RoBERTa and BART. These models require an environment with access to GPUs in order to effectively and efficiently train these models. All experiments were run on Google Colab.

### A. Hardware Specifications

All model training and evaluation were carried out in Google Colab using its GPU-supported runtime. The hardware configuration available during experimentation was:

- **Processing Unit (GPU):** NVIDIA Tesla T4 (16 GB VRAM) provided by Google Colab.
- **RAM:** 12–16 GB system memory.
- **Storage:** 50–80 GB temporary cloud storage allocated within Colab.
- **Execution Platform:** Entire code executed end-to-end in Google Colab environment.

This configuration was sufficient for fine-tuning transformer models for both phrase extraction and sentiment classification.

### B. Software Tools and Frameworks

The following software tools, libraries, and frameworks were utilized during implementation as shown in V:

HuggingFace Transformers and Datasets libraries were essential for loading pretrained models, tokenization, fine-tuning, and evaluation. PyTorch served as the backend framework for training.

| Software / Library | Version Used |
|---|---|
| Python | 3.10+ |
| PyTorch | 2.1.0+ |
| Transformers (HuggingFace) | 4.36+ |
| Datasets (HuggingFace) | 2.16+ |
| NumPy | 1.26+ |
| Pandas | 2.1+ |
| Matplotlib | 3.8+ |
| Scikit-learn | 1.3+ |
| Google Colab Runtime | GPU-enabled (T4 / L4) |

**TABLE V:** Software Tools and Libraries Used

### C. Development Environment

The full project development, training, and evaluation were performed using:

- **Primary Development Platform:** Google Colab (GPU Runtime)
- **IDE for Code Editing:** Visual Studio Code (optional)

Google Colab provided the required GPU acceleration, ease of library installation, and efficient runtime for transformer-based model training.

## X. CONSLUSION

This project successfully developed a two-stage NLP framework for fine-grained opinion mining by combining phrase extraction with sentiment classification. The first stage focused on identifying meaningful opinion-bearing phrases using transformer-based NER models such as RoBERTa, XLM-RoBERTa, BART, and BERT. Through 6-fold cross-validation and careful tuning to address overfitting, the system achieved stable performance with an overall accuracy of 0.82% and macro and weighted F1-scores of 0.72 and 0.82 respectively. The second stage applied the extracted phrases to sentiment classification models, which performed exceptionally well, achieving around 99% accuracy and a 99% F1-score, demonstrating the effectiveness of separating phrase extraction from sentiment prediction.

Overall, the proposed two-stage architecture proved to be more robust and interpretable than conventional sentence-level sentiment analysis. By isolating the exact opinion-bearing segments before classifying their polarity, the system provides more detailed insights and improves predictive reliability. The results confirm that transformer-based models, when combined with structured preprocessing and cross-validated training, form a powerful solution for fine-grained sentiment understanding in user-generated text.The results of this study reinforce the idea that transformer-based token classification and sentiment models are now mature enough to be deployed in real-world opinion mining systems.

## XI. SCOPES OF FUTURE WORK

Although the proposed pipeline achieves strong results, several directions for future improvement remain:

1) **Handling Token-Level Annotation Challenges:** critical to the effectiveness of the phrase extraction stage. When handled incorrectly, the creation of subword prefixes can introduce noise and misalignment between the token

and label. Future work should include extending the implementation of encoder-decoder or span-based natural language processing techniques to the processing of spans instead of just individual tokens.

2) **Expanding and Refining the Dataset:** The existing dataset is relatively limited and manually annotated, which resulted in some label skewing potentially affecting early training consistency, therefore creating a larger dataset with broader and more realistic reviews would aid in enhancing generalization. Methods for Semi-supervised learning can also be utilized to increase the training corpus with little manual input.

3) **Expanding and Refining the Dataset:** The present method is capable of producing accurate results when applied to much of the clean data collected for training, however there may be substantial variation in the real-world application of user reviews, which could include differing vocabulary, tone and complexity of sentences. Utilizing techniques for Domain Adaptation as a means for reducing the distribution shift of data between the training and deployment phases (e.g., fine-tuning on data collected from a specific domain, or using a Domain Adaptive Pre-training approach) will allow for the continued accuracy of the Text Extraction (Phrase Extraction) and Sentiment Classifier systems when applied to different/unknown domains.

4) **Methods to Incorporate Explainable AI:** Due to the inherent structure of transformer-based models, these models demonstrate excellent performance, however they are black box type models making it difficult to project or derive how a particular phrase or sentiment has been predicted. By integrating various Explainable AI methods such as Attention based Visualization methods, Layer-wise Relevance Propagation methods, or SHAP based Interpretability methods, users will be able to gain insight into the reasoning behind a transformer models prediction process. Such types of information will be especially critical to systems/solutions that require transparency and trust.

5) **Achieving Lower Computational Complexity for the Deployment Phase:** While transformer style models will generally yield a high degree of accuracy, the average computational resource requirement to undertake the inference operation on these models will also continue to be large in magnitude in order to perform the same level of service in a real-time manner or on a resource/location restricted basis. Future Work should also consider researching and identifying methods for Compression of transformer style models through Pruning, Quantization and Distilling methods to achieve more lightweight versions of the transformer models which are suitable for deployment onto Edge Device type configurations or restricted resource type configurations which will greatly improve the ability to scale this solution down for use in/through thousands of Off-the-Shelf Device or On-Device Applications.

In summary, while the developed system achieves strong quantitative performance and successfully validates the two-stage approach, several technical enhancements can further improve robustness, scalability, and real-world applicability.

## REFERENCES

[1] R. Sarawagi, "Information Extraction," *Foundations and Trends in Databases*, vol. 1, no. 3, pp. 261–377, 2008.

[2] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *Proc. NAACL-HLT*, 2019.

[3] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention Is All You Need," in *Proc. NeurIPS*, 2017.

[4] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "RoBERTa: A Robustly Optimized BERT Pretraining Approach," *arXiv preprint arXiv:1907.11692*, 2019.

[5] B. Liu, *Sentiment Analysis and Opinion Mining*. Synthesis Lectures on Human Language Technologies, 2012.

[6] J. Lafferty, A. McCallum, and F. Pereira, "Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data," in Proc. ICML, 2001.

[7] Z. Huang, W. Xu, and K. Yu, "Bidirectional LSTM-CRF Models for Sequence Tagging," arXiv:1508.01991, 2015.

[8] Y. Kim, "Convolutional Neural Networks for Sentence Classification," in Proc. EMNLP, 2014.

[9] Y. Liu et al., "RoBERTa: A Robustly Optimized BERT Pretraining Approach," arXiv:1907.11692, 2019.

[10] E. F. Tjong Kim Sang and S. Buchholz, "Introduction to the CoNLL-2000 Shared Task: Chunking," in Proc. CoNLL, 2000.

[11] B. Pang and L. Lee, "Opinion Mining and Sentiment Analysis," Foundations and Trends in Information Retrieval, vol. 2, no. 1–2, pp. 1–135, 2008.

[12] E. F. Tjong Kim Sang and F. De Meulder, "Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition," in *Proceedings of the Seventh Conference on Natural Language Learning*, 2003.

[13] N. Reimers and I. Gurevych, "Reporting Score Distributions Makes a Difference: Performance Study of LSTM-CRF Models for Sequence Tagging," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2017.

[14] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts, "Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank," in *Proceedings of EMNLP*, 2013.

[15] R. Kohavi, "A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection," in *Proceedings of the 14th International Joint Conference on Artificial Intelligence*, 1995.

[16] R. Kohavi, "A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection," in *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 1137–1143, 1995.

[17] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, and A. Rush, "Transformers: State-of-the-Art Natural Language Processing," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 38–45, 2020.

[18] **Souza, F., Nogueira, R., & Lotufo, R.**, BERTimbau: Pretrained BERT Models for Brazilian Portuguese, *Brazilian Conference on Intelligent Systems*, 2020.

[19] **Wang, B., & Lu, W.**, Learning Latent Structures for Semantic Parsing with Neural Networks, *AAAI*, 2018.

[20] **Zhang, M., & Liu, B.**, Aspect and Opinion Mining, *Springer*, 2018.

[21] **Sun, C., Huang, L., & Qiu, X.**, Utilizing BERT for Aspect-Based Sentiment Analysis via Constructing Auxiliary Sentence, *NAACL Workshop*, 2019.

[22] **Ribeiro, M. T., Singh, S., & Guestrin, C.**, Why Should I Trust You? Explaining the Predictions of Any Classifier, *KDD*, 2016.

[23] **Jacovi, A., & Goldberg, Y.**, Towards Faithfully Interpretable NLP Systems: How Should We Define and Evaluate Faithfulness?, *ACL*, 2020.

[24] **Xu, R., Yang, Y., & Lin, H.**, Deep Learning for Low-Resource Natural Language Processing, *IEEE Access*, 2020.

[25] **Gururangan, S., Marasović, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., & Smith, N. A.**, Don't Stop Pretraining: Adapt Language Models to Domains and Tasks, *ACL*, 2020.

[26] **Li, X., Wu, X., & Meng, L.**, Adversarial Learning for Robust Natural Language Understanding, *IEEE Transactions on Neural Networks and Learning Systems*, 2021.