

## Approach and Implementation

The task involved performing sentiment analysis on a dataset using the MARBERT model for feature extraction and fine-tuning. The dataset was split into training, validation, and test sets. The model was fine-tuned using the Hugging Face, and various evaluation metrics were computed to assess performance with **79.6% accuracy on test set**. Key steps:

- 1. Label Encoding:**
  - The sentiment labels were encoded using `LabelEncoder` to convert categorical labels into numerical format.
- 2. Data Splitting:**
  - The dataset was divided into training (80%), validation (10%), and test (10%) sets.
- 3. Model Selection:**
  - MARBERT was chosen for its effectiveness in handling Arabic text.
- 4. Tokenization:**
  - Text data was tokenized using MARBERT's tokenizer, ensuring consistency with the model's expected input format.
- 5. Dataset Creation:**
  - Custom PyTorch datasets were created for training, validation, and test sets.
- 6. Model Fine-Tuning:**
  - The model was fine-tuned with specific training arguments, including learning rate, batch size, and evaluation strategy.
- 7. Evaluation:**
  - Performance was evaluated on validation and test sets using accuracy, precision, recall, F1 score, and confusion matrix.

## Challenges Faced and Solutions

- 1. Handling Imbalanced Data:**
  - The dataset had an imbalanced distribution of labels, which could skew the model's performance. This was addressed by using stratified sampling during data splitting to maintain the label proportions in each set.
- 2. Tokenization Efficiency:**
  - Tokenizing large text data efficiently was crucial. The solution involved using batch tokenization provided by Hugging Face's tokenizer, which improved processing speed.
- 3. Model Selection and Compatibility:**
  - Ensuring compatibility between the tokenizer and model (MARBERT) was a challenge. Using the `AutoTokenizer` and `AutoModelForSequenceClassification` classes from Hugging Face ensured that the correct tokenizer-model pair was used.
- 4. Performance Optimization:**
  - Fine-tuning the model required careful selection of hyperparameters. Using the `TrainingArguments` class allowed for flexible configuration and optimization of training parameters.

## **Aspect-Based Sentiment Analysis Approach**

To adapt the current approach for aspect-based sentiment analysis, several changes would be necessary. The dataset must be annotated with aspect categories alongside sentiment labels, with each text sample having multiple labels for different aspects. The model architecture needs to be adapted for multi-label classification, treating each aspect category as a separate classification task. While tokenization would remain largely the same, it would include both text and aspect labels. A custom dataset class would be required to handle multi-label outputs, ensuring each sample has the appropriate aspect and sentiment labels. The training loop would be modified for multi-label classification, with evaluation metrics including aspect-wise precision, recall, and F1 scores.