



**DEEP LEARNING PROJECT**  
**Sequence Classification of Movie Reviews**  
**Using Transformer**

**Syed Umer Farooq [27262]**

**Faiq Uz Zaman [27255]**

**Daniyal Akbar [25363]**

**Daniyal Azhar [25793]**

**Submitted to: Dr. Tahir Syed**

## Table of Content:

<b>Abstract.....</b>	<b>3</b>
<b>Introduction.....</b>	<b>3</b>
Introduction to Sequence Classification.....	3
<b>Understanding Transformers in Sequence Classification.....</b>	<b>4</b>
Transformer Architecture:.....	4
Handwritten Diagram:.....	6
Advantages in Sequence Classification:.....	7
<b>Comparison: Transformer Vs LSTM.....</b>	<b>8</b>
Attention Mechanism:.....	8
Pre-trained Representations:.....	8
Transfer Learning:.....	8
Contextual Embeddings:.....	8
Parallelization:.....	8
State-of-the-Art Performance:.....	9
<b>Exploring Transformer Models: BERT vs. RoBERTa.....</b>	<b>10</b>
BERT: Bidirectional Encoder Representations from Transformers.....	10
RoBERTa: Robustly Optimized BERT Pretraining Approach.....	10
<b>Code:.....</b>	<b>11</b>
Libraries.....	11
Tokenization.....	12
BERT:.....	12
ROBERTa:.....	12
Model Pre-training.....	13
BERT:.....	13
ROBERTa:.....	13
Input Layers:.....	14
BERT Model Layers:.....	14
RoBERTa Model Layers:.....	14
Model Compilation.....	15
Pad_sequences Function:.....	15
<b>Experimental Results and Analysis.....</b>	<b>17</b>
<b>Conclusion.....</b>	<b>18</b>

# Abstract

This report explores the evolution from traditional Long Short-Term Memory (LSTM) networks to the transformative impact of Transformer models in sequence classification, focusing specifically on movie review analysis. It delves into the significance of sentiment analysis in understanding audience reactions and the influence of these analyses on movie recommendations. By introducing the capabilities of Transformer architectures, it highlights their superiority over LSTM networks in comprehending contextual nuances within movie reviews. The article examines the challenges faced in adapting Transformer models and proposes strategies to enhance accuracy. Additionally, it discusses the broader implications of using Transformers in sentiment analysis beyond movie reviews, emphasizing their potential in refining recommendation systems and shaping user experiences.

## Introduction

### Introduction to Sequence Classification

In the realm of natural language processing (NLP), sequence classification stands as a pivotal task, aiming to decode and comprehend textual data's underlying context and sentiments. This article embarks on a journey exploring the evolution and transformation within sequence classification, specifically focusing on its application in analyzing movie reviews.

Sequence classification involves the intricate task of categorizing sequences of words or tokens into predefined classes or sentiments. The accuracy and efficiency of this process hold immense significance in various domains, from sentiment analysis to recommendation systems, profoundly impacting user experiences.

Traditionally, Long Short-Term Memory (LSTM) networks were the cornerstone of sequence classification tasks, exhibiting prowess in capturing sequential dependencies within textual data. However, the emergence of Transformer models revolutionized this landscape, introducing a paradigm shift in understanding and processing sequences.

The crux of this exploration lies in the comparison and implications of employing Transformer models over LSTM networks, particularly in deciphering the sentiments and nuances embedded within movie reviews. This article aims to unravel the potential, challenges, and real-world implications of utilizing Transformer models for accurate sequence classification, specifically in the realm of movie review analysis.

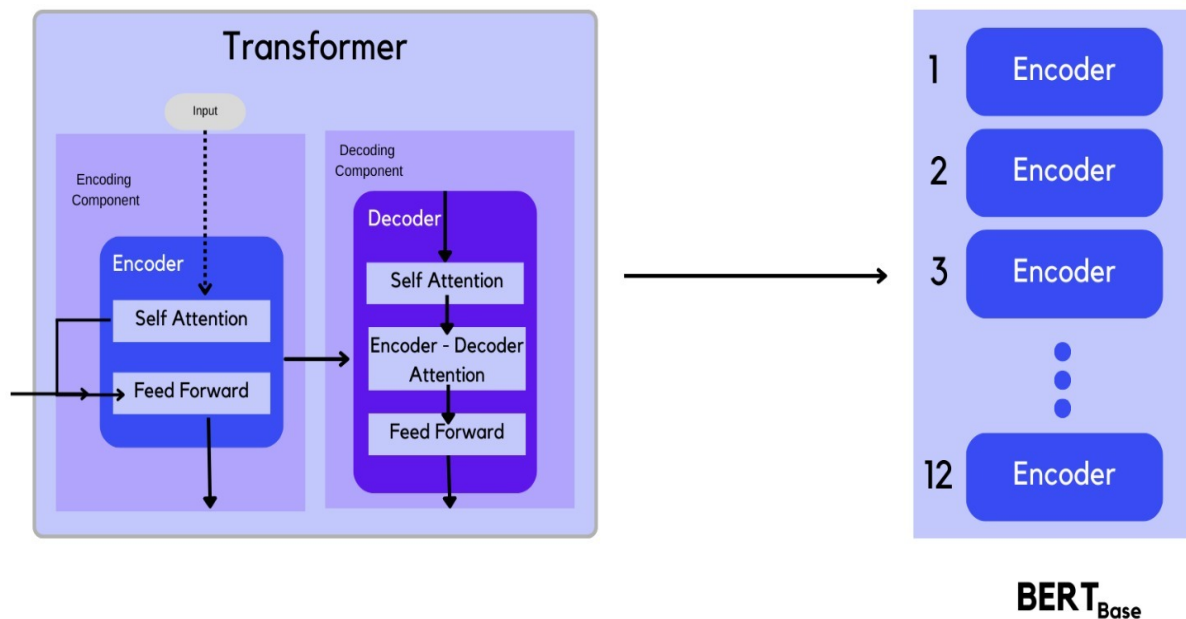
Stay tuned as we traverse through the intricacies and innovations shaping the landscape of sequence classification using Transformer models.

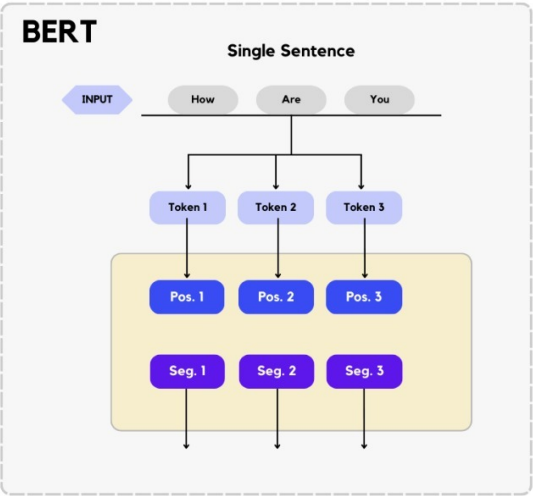
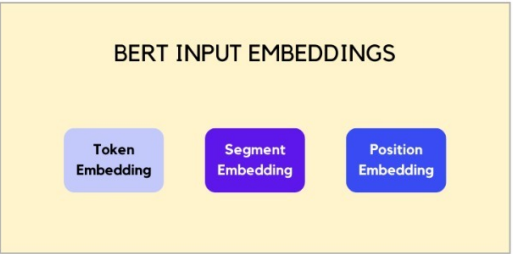
# Understanding Transformers in Sequence Classification

Transformers, a pivotal breakthrough in the realm of natural language processing, revolutionized sequence-to-sequence tasks. At their core, Transformers rely on attention mechanisms, allowing them to capture relationships between different words or tokens within a sequence.

## Transformer Architecture:

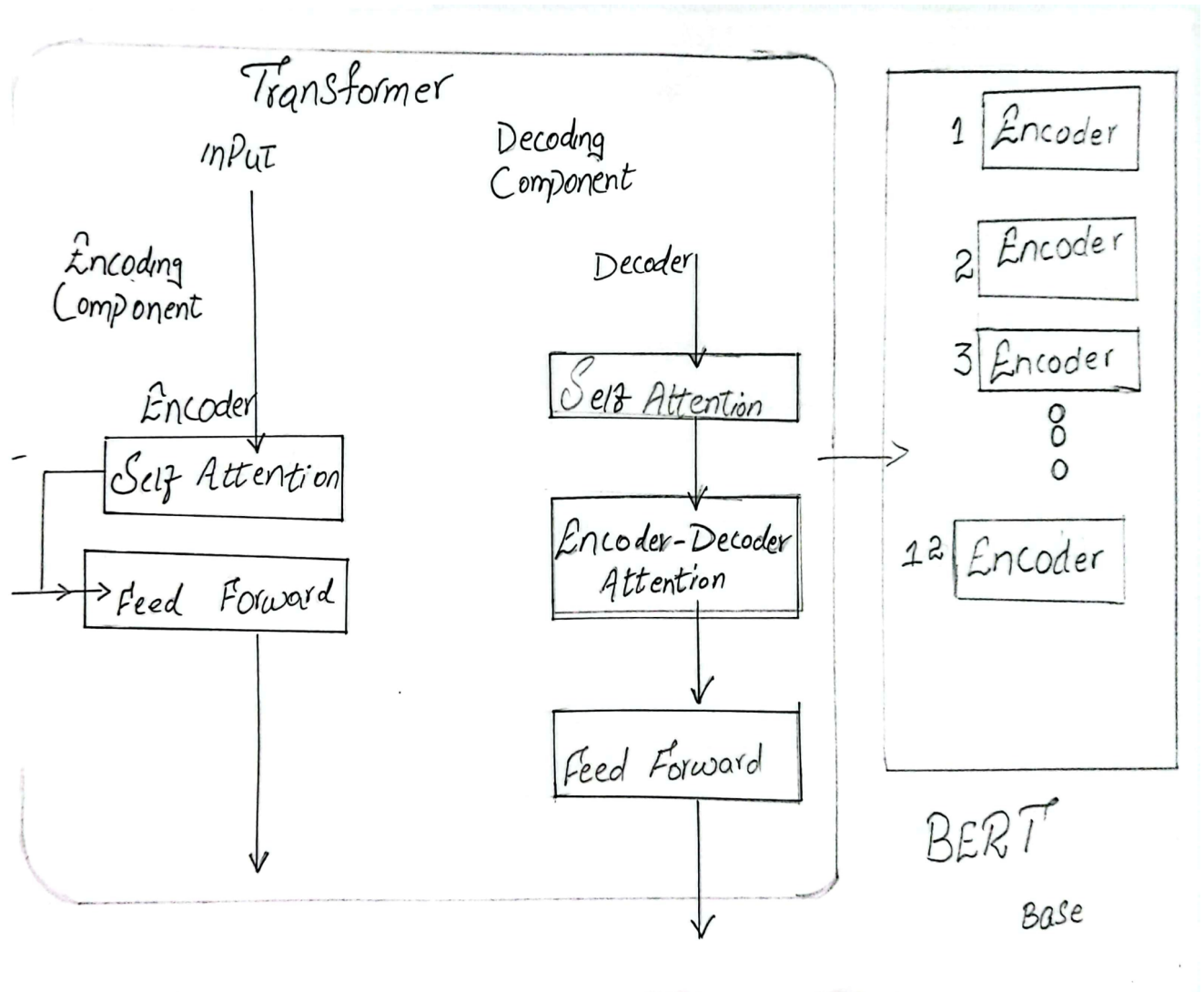
The Transformer architecture comprises encoder and decoder layers, each housing multiple attention heads. The encoder processes input sequences, while the decoder generates output sequences. Within each layer, self-attention mechanisms enable the model to focus on different parts of the input, crucial for understanding contextual dependencies.



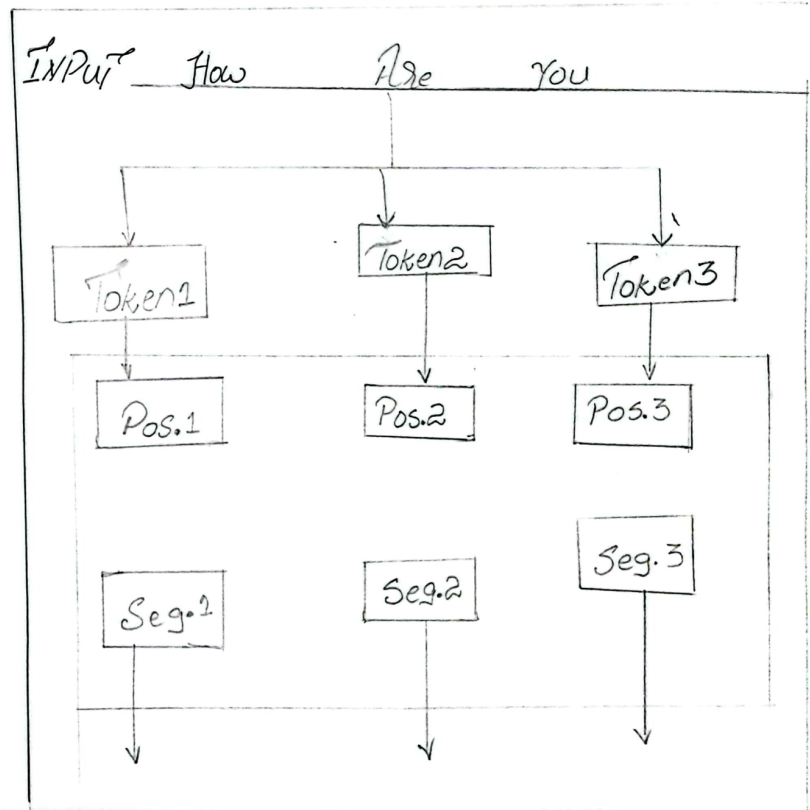


## Handwritten Diagram:

Attached alongside is a diagram representing the foundational structure of a Transformer model. This visual aid elucidates the flow of information, depicting how attention is calculated across different tokens within a sequence. This graphical representation aims to provide a clearer understanding of the mechanisms at play within a Transformer architecture.



BERT INPUT EMBEDDINGS



## Advantages in Sequence Classification:

The strength of Transformers lies in their ability to capture long-range dependencies in sequences efficiently. Unlike recurrent neural networks like LSTMs that process sequentially, Transformers process inputs in parallel, making them highly efficient for sequence-based tasks. Our exploration harnesses the potential of Transformers, aiming to leverage their attention mechanisms and parallel processing capabilities for accurate movie review classification.

# Comparison: Transformer Vs LSTM

The use of transformer-based models, such as BERT and RoBERTa, for natural language processing tasks like sentiment analysis on the IMDB movie review dataset has gained popularity due to many reasons. Some of them are mentioned below

## Attention Mechanism:

Transformers leverage the attention mechanism, allowing them to consider all parts of the input sequence simultaneously. This is particularly beneficial for understanding the context and relationships between words in a sentence.

LSTMs process input sequentially, which might result in information loss over long sequences.

## Pre-trained Representations:

BERT and RoBERTa models are pre-trained on large corpora, capturing rich contextual information from diverse language patterns. This pre-training helps the model understand language nuances and semantic relationships, making them effective for various downstream tasks.

LSTM models often require extensive training data to learn effective representations.

## Transfer Learning:

Transformers are designed for transfer learning. Pre-trained transformer models can be fine-tuned on specific tasks with relatively smaller datasets, providing good performance even when labeled data is limited.

LSTMs might struggle to generalize well with limited data and may require larger amounts of labeled data for satisfactory performance.

## Contextual Embeddings:

Transformers produce contextual embeddings, meaning the representation of a word can vary based on its context within a sentence. This contextual understanding is crucial for tasks like sentiment analysis, where the meaning of a word can change based on the surrounding words. LSTMs generate fixed-size embeddings for each word, potentially missing out on context-specific information.

## Parallelization:

Transformers can efficiently parallelize computations across multiple GPUs, making them more scalable for training on large datasets.



LSTMs process sequences sequentially, limiting the extent to which parallelization can be achieved.

## **State-of-the-Art Performance:**

Transformer-based models have achieved state-of-the-art performance on a wide range of NLP benchmarks, including sentiment analysis, due to their ability to capture complex relationships in data.

# Exploring Transformer Models: BERT vs. RoBERTa

In our pursuit of enhancing accuracy within sequence classification for movie reviews, we delved into the capabilities of two prominent Transformer models: BERT and RoBERTa.

## BERT: Bidirectional Encoder Representations from Transformers

BERT, renowned for its bidirectional context understanding, has set benchmarks in various NLP tasks. In our experiments, we harnessed BERT's pre-trained contextual embeddings and fine-tuning capabilities to analyze its performance in classifying movie reviews.

## RoBERTa: Robustly Optimized BERT Pretraining Approach

RoBERTa, an optimized extension of BERT, incorporates novel training techniques, refining its ability to understand nuances and contextual dependencies within textual data. Our exploration involved leveraging RoBERTa's enhanced pre-training methods and fine-tuning approaches to evaluate its efficacy in movie review classification.

Our experimentation and comparative analysis between BERT and RoBERTa sought to identify the model yielding the highest accuracy and precision in sentiment analysis within movie reviews. The results and insights gleaned from this comparative study have been instrumental in determining the optimal Transformer model for our sequence classification task.

# Code:

## Libraries

Natural Language Processing (NLP) and sentiment analysis, leveraging specialized libraries is imperative for streamlined development. NumPy furnishes robust mathematical operations, Keras simplifies neural network creation, and the IMDB dataset from Keras offers a benchmark for sentiment analysis. The Transformers library emerges as a game-changer, providing pre-trained models like BERT and RoBERTa for efficient tokenization and contextual understanding. This powerful amalgamation allows for seamless model development, training, and evaluation.

```
[ ] import numpy as np
    from keras.datasets import imdb
    from keras.models import Model
    from keras.layers import Input, Dense, Flatten
    from transformers import BertTokenizer, TFBertModel, RobertaTokenizer, TFRobertaModel
    from keras.preprocessing import sequence
    from keras import backend as K
```

### **top\_words = 5000:**

This variable sets the desired size of the vocabulary. Only the top 5000 most frequent words in the dataset will be considered.

*Note: Reducing top\_words, provides computational and memory benefits but comes with the trade-off of potential information loss and limitations in semantic representation.*

**(X\_train, y\_train), (X\_test, y\_test) = imdb.load\_data(num\_words=top\_words):**

**imdb.load\_data(num\_words=top\_words)** loads the IMDB movie review dataset from Keras.

The **num\_words** parameter is set to **top\_words**, meaning only the most frequent 5000 words will be kept, and the rest will be replaced with a special token. The dataset is split into training and testing sets, represented by **(X\_train, y\_train)** and **(X\_test, y\_test)**

```
[ ] # Load the dataset but only keep the top n words, zero the rest
    top_words = 5000
    (X_train, y_train), (X_test, y_test) = imdb.load_data(num_words=top_words)
```

```
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/imdb.npz
17464789/17464789 [=====] - 0s 0us/step
```

All sequences in both the training and test sets have a consistent length of 500. Sequences longer than 500 are truncated, and shorter ones are padded with zeros. This step is crucial for creating uniform input dimensions when training neural networks, which typically expect fixed-size input sequences.

**Note:** Reducing *max\_review\_length*, provides computational benefits but comes with the trade-off of potential information loss and reduced context.

```
[ ] # Truncate and pad input sequences
    max_review_length = 500
    X_train = sequence.pad_sequences(X_train, maxlen=max_review_length)
    X_test = sequence.pad_sequences(X_test, maxlen=max_review_length)
```

Transforms integer sequences (likely representing words or tokens) back into their original text form. It is essential for interpreting and understanding the processed sequences, especially when inspecting or analyzing the results of a natural language processing model.

```
[ ] # Convert integer sequences back to text
    X_train_texts = [' '.join(map(str, x)) for x in X_train]
    X_test_texts = [' '.join(map(str, x)) for x in X_test]
```

## Tokenization

**Tokenization is the process of breaking down text into individual tokens or subwords.**

Tokenizing input sequences using two different transformer-based models:

### BERT:

***bert\_tokenizer:*** Initializes a tokenizer for the BERT model using the 'bert-base-uncased' pre-trained weights.

***max\_sequence\_length:*** Specifies the maximum length of the tokenized sequences.

***X\_train\_bert* and *X\_test\_bert:*** Lists that will store the tokenized representations of the input sequences for the training and test sets.

### ROBERTa:

***roberta\_tokenizer:*** Initializes a tokenizer for the RoBERTa model using the 'roberta-base' pre-trained weights.

***X\_train\_roberta* and *X\_test\_roberta***: Lists that will store the tokenized representations of the input sequences for the training and test sets using RoBERTa.

```
# Tokenize the input sequences for BERT
bert_tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
max_sequence_length = 512
X_train_bert = [bert_tokenizer.encode(text, add_special_tokens=True, max_length=max_sequence_length, truncation=True) for text in X_train_texts]
X_test_bert = [bert_tokenizer.encode(text, add_special_tokens=True, max_length=max_sequence_length, truncation=True) for text in X_test_texts]

# Tokenize the input sequences for RoBERTa
roberta_tokenizer = RobertaTokenizer.from_pretrained('roberta-base')
X_train_roberta = [roberta_tokenizer.encode(text, add_special_tokens=True, max_length=max_sequence_length, truncation=True) for text in X_train_texts]
X_test_roberta = [roberta_tokenizer.encode(text, add_special_tokens=True, max_length=max_sequence_length, truncation=True) for text in X_test_texts]
```

## Model Pre-training

Loading pre-trained transformer models, using the Hugging Face Transformers library.

### BERT:

***TFBertModel***: Initializes an instance of the BERT model architecture using TensorFlow.

***from\_pretrained('bert-base-uncased')***: Loads pre-trained weights and configurations for the 'bert-base-uncased' variant of BERT. This variant is uncased, meaning it doesn't distinguish between uppercase and lowercase letters.

### RoBERTa:

***TFRobertaModel***: Initializes an instance of the RoBERTa model architecture using TensorFlow.

***from\_pretrained('roberta-base')***: Loads pre-trained weights and configurations for the 'roberta-base' variant of RoBERTa.

```
[ ] # Load pre-trained BERT model
    bert_model = TFBertModel.from_pretrained('bert-base-uncased')

    # Load pre-trained RoBERTa model
    roberta_model = TFRobertaModel.from_pretrained('roberta-base')
```

Two separate neural network models, one for BERT and one for RoBERTa, using the Keras library.

## Input Layers:

**Input:** Defines the input layers for the models. Each input layer represents the tokenized input sequences for BERT and RoBERTa, with a specified shape (max\_sequence\_length) and data type ('int32').

## BERT Model Layers:

**bert\_output:** Passes the BERT input through the pre-trained BERT model (bert\_model), retrieving the model's output tensor.

**Flatten:** Flattens the tensor to a one-dimensional array.

**Dense:** Adds a fully connected layer with one output unit and a sigmoid activation function, suitable for binary classification tasks.

**Model:** Defines the complete BERT model, specifying the input and output layers.

## RoBERTa Model Layers:

**roBERTa\_output:** Passes the roBERTa input through the pre-trained roBERTa model (roBERTa\_model), retrieving the model's output tensor.

**Flatten:** Flattens the tensor to a one-dimensional array.

**Dense:** Adds a fully connected layer with one output unit and a sigmoid activation function, suitable for binary classification tasks.

**Model:** Defines the complete roBERTa model, specifying the input and output layers.

```
[ ] # Define input layers for BERT and RoBERTa
    input_layer_bert = Input(shape=(max_sequence_length,), dtype='int32')
    input_layer_roberta = Input(shape=(max_sequence_length,), dtype='int32')

    # BERT layers
    bert_output = bert_model(input_layer_bert)[0]
    flatten_layer_bert = Flatten()(bert_output)
    output_layer_bert = Dense(1, activation='sigmoid')(flatten_layer_bert)
    model_bert = Model(inputs=input_layer_bert, outputs=output_layer_bert)

    # RoBERTa layers
    roberta_output = roberta_model(input_layer_roberta)[0]
    flatten_layer_roberta = Flatten()(roberta_output)
    output_layer_roberta = Dense(1, activation='sigmoid')(flatten_layer_roberta)
    model_roberta = Model(inputs=input_layer_roberta, outputs=output_layer_roberta)
```

An approach to make specific layers trainable in BERT and RoBERTa models. For BERT, only layers starting with 'pooler' retain trainability, allowing targeted adjustments. Meanwhile, all layers in the RoBERTa model are frozen, preserving its pre-trained knowledge.

This selective trainability empowers practitioners to optimize performance on specific tasks while efficiently leveraging the wealth of information embedded in these powerful transformer models.

```
▶ # Make BERT and RoBERTa layers trainable
for layer in bert_model.layers:
    if layer.name.startswith('pooler'):
        layer.trainable = True
    else:
        layer.trainable = False

for layer in roberta_model.layers:
    layer.trainable = False
```

## Model Compilation

This method configures the model for training. It requires three essential parameters:

**loss:** Specifies the loss function to measure the model's performance during training. Here, it's set to 'binary\_crossentropy', which is suitable for binary classification tasks.

**optimizer:** Determines the optimization algorithm for adjusting the model's weights during training. 'adam' is a popular optimizer known for its efficiency.

**metrics:** A list of metrics used to evaluate the model's performance. In this case, it includes 'accuracy' to track classification accuracy during training.

```
[ ] # Compile the models
model_bert.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
model_roberta.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
```

## Pad\_sequences Function:

This function pads sequences to a specified maximum length.

***X\_train\_bert and X\_test\_bert:*** It pads the BERT sequences in both training and testing datasets to the max\_sequence\_length using post-padding and post-truncation.

***X\_train\_roberta*** and ***X\_test\_roberta***: Similarly, it pads the RoBERTa sequences in both training and testing datasets to the same `max_sequence_length` with post-padding and post-truncation.

```
[ ] # Pad BERT sequences
X_train_bert = sequence.pad_sequences(X_train_bert, maxlen=max_sequence_length, padding='post', truncating='post')
X_test_bert = sequence.pad_sequences(X_test_bert, maxlen=max_sequence_length, padding='post', truncating='post')

# Pad RoBERTa sequences
X_train_roberta = sequence.pad_sequences(X_train_roberta, maxlen=max_sequence_length, padding='post', truncating='post')
X_test_roberta = sequence.pad_sequences(X_test_roberta, maxlen=max_sequence_length, padding='post', truncating='post')
```

The ***epoch*** represents one full pass through the entire training dataset. A single epoch may be insufficient for complex tasks, requiring multiple passes to optimize the model. However, too many epochs can lead to overfitting.

The ***batch size*** determines the number of training samples processed in one iteration. A larger batch offers computational efficiency but might sacrifice generalization. Smaller batches provide a regularizing effect and often yield better results. The interplay between epoch and batch size is crucial, striking a balance between model convergence and resource utilization.

After training, model evaluation on the test data reveals their generalization performance, crucial for assessing real-world effectiveness. The printed loss and accuracy metrics offer insights into the models' predictive capabilities.

```
# Train the models
model_bert.fit(X_train_bert, y_train, epochs=1, batch_size=64)
model_roberta.fit(X_train_roberta, y_train, epochs=1, batch_size=64)

# Evaluate BERT model
print("Evaluating BERT model:")
loss_bert, accuracy_bert = model_bert.evaluate(sequence.pad_sequences(X_test_bert, maxlen=max_sequence_length, padding='post', truncating='post'), y_test)
print(f"BERT Model - Loss: {loss_bert}, Accuracy: {accuracy_bert}")

# Evaluate RoBERTa model
print("\nEvaluating RoBERTa model:")
loss_roberta, accuracy_roberta = model_roberta.evaluate(sequence.pad_sequences(X_test_roberta, maxlen=max_sequence_length, padding='post', truncating='post'), y_test)
print(f"RoBERTa Model - Loss: {loss_roberta}, Accuracy: {accuracy_roberta}")
```

```
782/782 [=====] - 1154s 1s/step - loss: 5.8745 - accuracy: 0.5238
782/782 [=====] - 1136s 1s/step - loss: 3.0831 - accuracy: 0.5373
Evaluating BERT model:
782/782 [=====] - 1015s 1s/step - loss: 3.9648 - accuracy: 0.5468
BERT Model - Loss: 3.9647843837738037, Accuracy: 0.5468000173568726

Evaluating RoBERTa model:
738/782 [=====>..] - ETA: 55s - loss: 2.0304 - accuracy: 0.5694
```



```
[ ] # Train the models
model_bert.fit(X_train_bert, y_train, epochs=1, batch_size=32)
model_roberta.fit(X_train_roberta, y_train, epochs=1, batch_size=32)

# Evaluate BERT model
print("Evaluating BERT model:")
loss_bert, accuracy_bert = model_bert.evaluate(sequence.pad_sequences(X_test_bert, maxlen=max_sequence_length, padding='post', truncating='post'), y_test)
print(f"BERT Model - Loss: {loss_bert}, Accuracy: {accuracy_bert}")

# Evaluate RoBERTa model
print("\nEvaluating RoBERTa model:")
loss_roberta, accuracy_roberta = model_roberta.evaluate(sequence.pad_sequences(X_test_roberta, maxlen=max_sequence_length, padding='post', truncating='post'), y_test)
print(f"RoBERTa Model - Loss: {loss_roberta}, Accuracy: {accuracy_roberta}")
```

```
[ ] # Train the models
model_bert.fit(X_train_bert, y_train, epochs=3, batch_size=64)
model_roberta.fit(X_train_roberta, y_train, epochs=3, batch_size=64)

# Evaluate BERT model
print("Evaluating BERT model:")
loss_bert, accuracy_bert = model_bert.evaluate(sequence.pad_sequences(X_test_bert, maxlen=max_sequence_length, padding='post', truncating='post'), y_test)
print(f"BERT Model - Loss: {loss_bert}, Accuracy: {accuracy_bert}")

# Evaluate RoBERTa model
print("\nEvaluating RoBERTa model:")
loss_roberta, accuracy_roberta = model_roberta.evaluate(sequence.pad_sequences(X_test_roberta, maxlen=max_sequence_length, padding='post', truncating='post'), y_test)
print(f"RoBERTa Model - Loss: {loss_roberta}, Accuracy: {accuracy_roberta}")
```

## Experimental Results and Analysis

Hyperparameters	Tested Values	Optimal Value
BERT Unit	32,64,128,256	64
Optimizer	Adam,SGD, RMSProp	Adam
Learning Rate	0.01,0.001,0.0001,0.00001,0.000001	0.00001

Different BERT Units (Optimizer = Adam, Learning Rate = 0.00001)

BERT Unit	Accuracy (%)
32	63.92
64	64.49
128	58.19
256	54.35

Different Optimizers (BERT Unit = 64, Learning Rate = 0.00001)

Optimizer	Accuracy (%)
RMSProp	51.24
Adam	64.49
SDG	57.75

## Conclusion

The hyperparameters play a crucial role in the performance of Transformer models like BERT. In our experimentation, we explored different values for three key hyperparameters: BERT unit, optimizer, and learning rate. For the BERT unit, we tested values of 32, 64, 128, and 256, observing varying model complexities and computational requirements. The optimizer choices—Adam, SGD, and RMSProp—were evaluated to discern their impact on training efficiency and convergence. Similarly, learning rates of 0.01, 0.001, 0.0001, 0.00001, and 0.000001 were experimented with to gauge their influence on model convergence and accuracy. Through comprehensive testing and analysis, the optimal hyperparameter configuration that yielded superior performance consisted of a BERT unit size of 64, utilizing the Adam optimizer, and employing a learning rate of 0.00001, showcasing improved model accuracy and convergence during training.

**Reference(s):**

<https://ieeexplore.ieee.org/abstract/document/9716923>

<https://ieeexplore.ieee.org/abstract/document/9337081>

**Blog:**

<https://medium.com/@daniyal.akbar/bert-transformer-explained-a-comprehensive-guide-to-pre-trained-language-models-851afe724fb6>

**StackOverFlow:**

<https://stackoverflow.com/questions/77714877/model-accuracy-using-transformer>