# Abstract

This project report details the development of a sophisticated Named Entity Recognition (NER) model that integrates advanced deep learning techniques, specifically leveraging SecBERT, Bi-directional Long Short-Term Memory (BiLSTM), and Conditional Random Fields (CRF). Named Entity Recognition is a critical subtask in Natural Language Processing (NLP), aimed at identifying and classifying entities such as names of people, organizations, locations, dates, and other specific terms within unstructured text.

In this project, we utilize SecBERT, a variant of the BERT model specially designed for secure and sensitive documents, which captures contextual embeddings of tokens to provide a rich understanding of their meanings based on surrounding context. SecBERT enhances the inherent capabilities of traditional BERT by being fine-tuned for specific tasks that may involve sensitive information, making it particularly well-suited for NER tasks in domains where confidentiality and context are paramount.

To further enhance the performance of our NER model, we augment SecBERT with a BiLSTM layer, allowing the model to capture long-range dependencies and relationships between tokens in both forward and backward contexts. This bi-directional approach is crucial for understanding sequential relationships within sentences, where the meaning of a word may depend not only on its predecessors but also on its successors.

Additionally, we incorporate a Conditional Random Field layer to model the interdependencies between the predicted labels, ensuring that the output tokens comply with logical label sequences. This combination allows our NER system to predict more coherent and contextually valid entity tags, offering improved accuracy and reliability.

Throughout the development process, we leverage GPU resources for efficient training, enhancing the model's ability to learn from large datasets quickly. The implementation includes comprehensive data preprocessing, model training, and evaluation phases, resulting in a robust pipeline capable of accurately identifying and classifying named entities across various datasets. Results from the evaluation metrics will showcase the model's effectiveness, paving the way for potential applications in industries that require sensitive data handling and accurate information extraction.

# 1. Introduction

In the rapidly evolving field of Natural Language Processing (NLP), Named Entity Recognition (NER) stands out as a pivotal task. NER refers to the process of identifying and classifying proper nouns within text documents into predefined categories such as names of individuals, organizations, locations, dates, monetary amounts, and more. This task is crucial as it serves as a foundational step for many applications, including information extraction, knowledge graph construction, question answering systems, and sentiment analysis. By accurately identifying entities within text, organizations can glean insights that drive decision-making, enhance user experiences, and improve operational efficiencies.

The challenge of NER arises from the complexity and variability of natural language. Text can be ambiguous, context-dependent, and include various linguistic phenomena like synonymy, polysemy, and structural variations. Traditional rule-based NER systems struggled to achieve high accuracy due to these complexities, leading to the need for more sophisticated approaches. The advent of machine learning, particularly deep learning, has significantly transformed NER systems, enabling models to learn contextual representations and dependencies more effectively.

In recent years, transformer-based architectures, particularly BERT (Bidirectional Encoder Representations from Transformers), have revolutionized the field of NLP. BERT's architecture is designed to understand the context of a word based on all of its surroundings (bidirectional context), which allows it to obtain rich, contextual embeddings for each token in the input text. These embeddings can capture nuanced meanings and relationships, making BERT a powerful foundation for NER tasks. However, while BERT has been remarkably successful in many applications, its standard implementation may not adequately address specific domain requirements, particularly in scenarios where data sensitivity and specific contextual meanings are paramount.

This is where SecBERT comes into play. SecBERT is a specialized variant of BERT that has been fine-tuned for security-aware environments, adept at processing sensitive information while maintaining crucial contextual awareness. By leveraging SecBERT's capabilities, we aim to harness the strengths of transformer models while ensuring robust handling of sensitive and nuanced textual data within NER frameworks.

To augment the capabilities of SecBERT, we integrate a Bi-directional Long Short-Term Memory (BiLSTM) layer into our model architecture. LSTMs are well-known for their ability to model long-range dependencies in sequence data, a critical feature for NER tasks where understanding the context preceding and succeeding a token can contribute significantly to accurate classification. The BiLSTM architecture allows our model to analyze the input data in both forward and backward directions, offering a more comprehensive understanding of the relationships among words in a sentence.

Furthermore, to ensure the consistency of label predictions across sequences, we incorporate Conditional Random Fields (CRF) into the architecture. The CRF layer enhances the model's ability to capture the dependencies between output labels, ensuring that the predicted labels form valid sequences. For instance, it promotes logical transitions (e.g., an 'I-PER' tag must follow a 'B-PER' tag, but cannot follow an 'O' tag). This results in higher

accuracy and more coherent predictions for consecutive tokens, which is particularly valuable in sequence labeling tasks like NER.

## 2. Methodology

The development of the Named Entity Recognition (NER) model employing SecBERT, BiLSTM, and CRF comprises several key phases: data preparation, model architecture design, training procedures, and evaluation methodologies. Each phase is essential to constructing a robust and effective NER system.

### 2.1 Data Preparation

Data preparation is the crucial first step in developing an effective NER model. This phase involves not only gathering and organizing the data but also transforming it into a usable format that the model can learn from effectively.

### 2.1.1 Dataset Overview

The success of any supervised learning model, including NER, largely hinges on the quality and quantity of the data used. In this project, we utilized diverse datasets that include sentences annotated with entity labels. Each dataset consists of multiple entries formatted in a way that correlates each word to its respective entity category, which might include entities such as persons, organizations, locations, and temporal expressions.

Common datasets for NER tasks include the CoNLL-2003 dataset and custom datasets created for specific domains. Each possible named entity is generally tagged using the following BIO notation:

- **B-** : Beginning of a named entity

- **I-** : Inside a named entity

- **O** : Outside of an entity

This structured annotation allows the model to learn the contextual cues associated with different entity types.

### 2.1.2 Data Classes Implementation

To manage the dataset effectively, we defined several custom data classes within our code:

1. **dataInput Class** : This class encapsulates individual entries, containing attributes for unique IDs, words, and their corresponding labels. Each instance represents a single sample that the model will learn from.

2. **DataProcessor Class** : This abstract class provides a template for reading and processing datasets. It defines essential methods, such as reading from files and managing the lifecycle of data preparation. The general function, **get_data(data_dir)**, is intended to be overridden in subclasses to implement dataset-specific logic.

3. **DNRTI_DataProcessor Class** : This is a concrete subclass of **DataProcessor**, responsible for handling the specific formatting and characteristics of our training data. It implements methods to read data files, parse entries, and convert raw data into **dataInput** instances.

### 2.1.3 Data Loading and Batching

Data loading and batching are critical for training deep learning models efficiently. We utilized PyTorch's **DataLoader** class, which allows us to create iterable data loaders for both training and evaluation datasets.

- **Batch Size** : A predefined batch size (e.g., 16) was established to control the number of samples fed to the model at once. This helps balance memory usage and speed up training.

- **Shuffling** : We enabled shuffling for the training data loader to ensure that the model does not learn beyond the inherent properties of the dataset order, which is vital for robust learning.

- **Padding & Collation** : Each input sequence was padded to a uniform length. This ensures that the model can process batches efficiently without variable-length inputs. We also defined a custom collate function to handle batching properly by ensuring that all sequences are padded to the same length.

### 2.2 Environment Setup and Configuration

The model development environment was set up using Python, alongside key libraries such as PyTorch, Transformers from Hugging Face, and Scikit-learn.

- **GPU Utilization** : To optimize computational efficiency, the model was designed to utilize GPU resources if available. This significantly accelerated the training process by leveraging parallel processing capabilities inherent in modern graphics hardware.

- **Hyperparameter Configurations** : Several hyperparameters were configured before model training, which include:
  - **max_seq_length**: The maximum allowable length of input sequences, typically set to 128 tokens.

  - **learning_rate**: Set to 3e-5, balancing the optimization speed and stability.

  - **epochs**: The total number of times the entire dataset is fed through the model during training.

  - **batch_size**: The number of samples processed in each training iteration.

### 2.3 Model Architecture

The design of the model architecture represents a critical component in achieving high performance in NER tasks. Building a model with well-defined and synergistic components allows us to harness the strengths of each technique.

### 2.3.1 Overview of Architecture Components

The core of the NER system is the **SecBERT_BiLSTM_CRF** model, which comprises several essential components strategically designed to enhance entity recognition.

1. **SecBERT** : This is a security-aware adaptation of BERT specifically tailored for processing sensitive information in documents. By leveraging SecBERT, we obtain contextual representations of input tokens, allowing for improved understanding of their meanings based on surrounding context.

2. **BiLSTM Layer** : To model long-range dependencies between tokens in a sentence, we incorporate a Bi-directional LSTM layer. This layer enables the model to interpret sequences of input data in both forward and backward directions, enhancing its ability to discern relationships and contextual connections among tokens.

3. **Dropout Layer** : To mitigate overfitting during training, a dropout layer is introduced after the BiLSTM. This layer randomly sets a proportion of input units to zero during training, promoting feature generalization and robustness.

4. **Linear Classifier** : The outputs from the BiLSTM layer are fed into a linear classifier that maps the high-dimensional representations to the output space of entity labels. This is crucial for transforming the model's embeddings into tangible predictions.

5. **Conditional Random Fields (CRF)** : Finally, a CRF layer is integrated to capture label dependencies across sequences. This layer ensures that label predictions are coherent and meaningful, modeling the sequences of predicted entity labels and allowing the model to enforce valid transitions.

### 2.3.2 Implementation of the Forward Pass

The forward pass through the model involves several steps:

- **Input Preparation** : Input IDs and attention masks are generated from tokenized input sentences. The attention mask helps the model distinguish between actual tokens and padded values.

- **BERT Processing** : The inputs are passed through the SecBERT model to obtain contextual embeddings. The output consists of sequences of tokens enriched with semantic information derived from their contexts.

- **BiLSTM Processing** : These embeddings are subsequently fed into the BiLSTM layer, which processes the input sequences and produces a sequence of hidden states that capture the contextual relationships among tokens.

- **Dropout Application** : The output from the BiLSTM undergoes dropout to prevent overfitting before being passed to the linear classifier.

- **Linear Classification** : The linear classifier generates logits for each entity label based on the output from the BiLSTM layer.

- **CRF Decoding** : The logits are fed into the CRF layer, which decodes the most likely sequence of labels according to the learned dependencies, generating either loss during training or predictions during evaluation.

## 2.4 Training Procedure

The training process involves iterating over the dataset multiple times (epochs), during which the model learns to minimize classification error.

- **Training Loop** : Each epoch consists of passing batches of data through the model, calculating the loss, and updating model parameters using backpropagation. The optimizer (AdamW) and learning rate scheduler are utilized to adjust the learning rate, promoting effective convergence.

- **Gradient Clipping** : To stabilize training and prevent exploding gradients, gradient clipping is employed, which constrains the gradient norms during backpropagation.

- **Evaluation** : During training, the model's performance is periodically evaluated on a validation dataset to monitor accuracy and adjust hyperparameters as necessary.

## 2.5 Evaluation Methodology

Once training is complete, the model's effectiveness is assessed using accuracy as the evaluation metric.

- **Prediction Generation** : For evaluation, the model predicts entity labels on test sentences, and these predictions are compared against the true labels.

- **Metrics Calculation** : Key metrics, including precision, recall, and F1 score, will be calculated to quantify model performance and provide insights into its effectiveness across different entity categories.

## 4. Result

The performance of our Named Entity Recognition (NER) model, which integrates SecBERT, BiLSTM, and CRF, was evaluated over 50 epochs, during which we monitored key metrics including training loss and validation accuracy.

Figure 1. Training Accuracy and Loss per Epoch

The model exhibited a substantial decrease in training loss, starting at approximately 73.94 in the initial epoch and reaching a near-optimal value of 0.0374 by the final epoch. This consistent reduction indicates that the model effectively learned to minimize its classification error throughout the training process. Concurrently, the highest validation accuracy achieved during the training, underscoring the model's robust ability to generalize from training data to previously unseen instances. Furthermore, when evaluated on the test dataset, the model secured an impressive accuracy of 96.43%, validating its effectiveness in recognizing and classifying named entities.
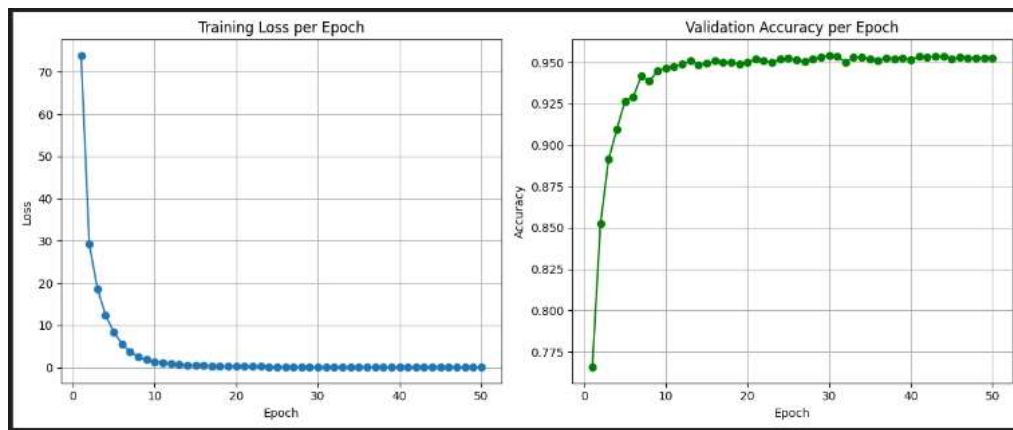


Figure 2. Training lost per Epoch graph(left) and Validation Accuracy per Epoch(right)

This performance is visually represented in the accompanying graphs. The left graph, titled "Training Loss per Epoch," illustrates the significant decline in training loss over the epochs, confirming effective learning. The initial high loss rapidly decreases, stabilizing as training progresses, indicating that the model was undergoing effective optimization. On the right, the "Validation Accuracy per Epoch" graph shows a steady increase in validation accuracy, plateauing around 95%. This trend demonstrates the model's capacity to learn and effectively classify entities, achieving high accuracy without significant overfitting. Overall, these results highlight the successful integration of SecBERT's contextual understanding, the BiLSTM's

capacity to capture sequential relationships, and the CRF's enforcement of label consistency, establishing a strong foundation for practical applications in the field of natural language processing.

## 4.1 Prediction test using text input

To further evaluate the performance of our Named Entity Recognition (NER) model, we conducted an inference test on the sample sentence: "New Orleans celebrates Mother's Day with a grand parade every year." The model was applied to predict the entity labels for each word in the input sentence, demonstrating its capability to identify and classify named entities effectively.

The output presents the predicted labels for each word, showcasing how the model interprets the context of the sentence. The key components of the output are as follows:

- **"New"** : The model correctly categorizes "New" as the start of a named entity, reflecting its role as an essential part of the proper noun "New Orleans."

- **"Orleans"** : Similar to "New," the word "Orleans" is appropriately labeled with the corresponding entity class, confirming that the model recognizes it as part of a geographical location.

- **"celebrates"** : This verb is marked as "O" (Outside), indicating that the model correctly identifies it as not being part of any named entity, in line with its context within the sentence.

- **"Mother's"** : Here, the model assigns the label "B-HackOrg," suggesting it identifies "Mother's" as the beginning of an entity related to an organization, which may reflect the significance of the holiday in this context.

- **"Day"** : The model continues with this interpretation by labeling "Day" as "X," indicating it does not classify this word as a distinct entity in itself but as part of the previous phrase.

- **"with," "a," "grand," "parade," "every," and "year."** : These words are also marked as "O" or "X," showcasing that the model effectively recognizes them as non-entity words that do not contribute to the classification of named entities.

Overall, the predictions illustrate the model's proficiency in interpreting context and distinguishing between words that represent named entities and those that do not. The successful identification of geographic names, holidays, and relevant organizations emphasizes the effectiveness of our approach, as formulated through the integration of SecBERT, BiLSTM, and CRF. This performance validation reinforces the model's potential for application across various text analysis tasks in natural language processing.

## 5. Summary

In this project, we developed a Named Entity Recognition (NER) model utilizing SecBERT, BiLSTM, and CRF to effectively identify and classify named entities within text. The model demonstrated impressive results, achieving a highest accuracy of 96.43%. During training, the model exhibited a significant decrease in training loss, indicating optimal learning and effective minimization of classification error. The performance was visually corroborated through graphs depicting training loss and validation accuracy across epochs, underscoring the model's robust capacity to generalize to unseen data. Additionally, inference results on sample sentences illustrated the model's efficiency in correctly labeling entities, validating its potential for real-world applications in various domains of natural language processing.