

# Invoice Data Extraction Report

## Table of Contents

1. Introduction
2. Approach
3. Model Architecture
4. Training Process
  - Data Collection
  - Data Preprocessing
  - Data Annotation
  - Fine-Tuning the Model
5. Evaluation Metrics
6. Optimization Techniques
  - Model Conversion to ONNX
  - Model Quantization
7. Deployment Steps
8. Performance
  - Inference Time
  - Accuracy
9. Conclusion

---

## 1. Introduction

This report documents the approach, model architecture, training process, evaluation metrics, optimization techniques, deployment steps, and performance of a machine learning model designed to extract key information from invoices. The solution aims to handle various invoice formats in English without hardcoded labels, leveraging the contextual understanding capabilities of a fine-tuned BERT model.

## 2. Approach

The project is divided into three main parts:

1. **Model Training:** Involves data collection, preprocessing, annotation, and fine-tuning a pre-trained BERT model.

2. **Model Optimization:** Focuses on converting the model to ONNX format and applying quantization techniques to improve efficiency.
3. **Model Deployment:** Ensures the optimized model runs effectively on a client desktop.

### 3. Model Architecture

We selected a pre-trained BERT model from Hugging Face's Transformers library (`dbmdz/bert-large-cased-finetuned-conll103-english`). BERT (Bidirectional Encoder Representations from Transformers) is a transformer-based model designed to understand the context of words in a sentence, making it ideal for tasks involving natural language understanding.

#### Key Features:

- **Bidirectional Context Understanding:** Captures context from both directions (left-to-right and right-to-left).
- **Fine-Tuning Capability:** Can be fine-tuned on specific tasks with relatively small datasets.
- **Transfer Learning:** Leverages pre-trained knowledge, reducing the need for extensive labeled data.

### 4. Training Process

#### Data Collection

A diverse dataset of invoices in PDF format was collected from the internet. The dataset included invoices from different companies and with varying formats to ensure the model could generalize well.

#### Data Preprocessing

1. **Text Extraction:** OCR (Optical Character Recognition) using `pdfplumber` and `pytesseract` was employed to convert PDFs to text.
2. **Text Cleaning:** The extracted text was cleaned to remove unnecessary whitespace, special characters, and line breaks.

#### Data Annotation

Manual annotation was performed to create a JSON file (`annotated_data.json`) containing key information such as sender, receiver, VAT number, and amounts. Each entry in the JSON file was structured as follows:

```
json
```

```
Copy code
```

```
{  
  
    "text": "Invoice from ABC Ltd. to XYZ Ltd. Total: $1000. VAT:  
123456789.",
```

```
    "annotations": {  
        "sender": "ABC Ltd.",  
        "receiver": "XYZ Ltd.",  
        "amount": "$1000",  
        "vat_number": "123456789"  
    }  
}
```

## Fine-Tuning the Model

The annotated data was used to fine-tune the pre-trained BERT model. The training parameters were set as follows:

- **Learning Rate:**  $2e-5$
- **Batch Size:** 16
- **Epochs:** 3

The dataset was split into training and validation sets to evaluate the model's performance during training.

## 5. Evaluation Metrics

The model's performance was evaluated using the following metrics:

- **Precision:** The ratio of correctly predicted positive observations to the total predicted positives.
- **Recall:** The ratio of correctly predicted positive observations to all observations in the actual class.
- **F1-Score:** The weighted average of Precision and Recall.

## 6. Optimization Techniques

### Model Conversion to ONNX

The fine-tuned model was converted to ONNX format for efficient deployment using the `transformers` library and the `onnx` package.

### Model Quantization

Dynamic quantization was applied to reduce the model size and improve inference speed using `onnxruntime`. This step ensures the model runs efficiently on client desktops without significant loss in accuracy.

## 7. Deployment Steps

1. **Set Up Client Environment:** Ensure the client machine has the necessary libraries installed (`onnx`, `onnxruntime`, `transformers`).

## 2. Load and Run the Model:

- Load the optimized ONNX model using ONNX Runtime.
- Run the model on sample invoice text to extract key information.

## 8. Performance

### Inference Time

The optimized model showed a reduction in inference time by approximately 50% compared to the original model, demonstrating the effectiveness of the optimization techniques.

### Accuracy

The accuracy remained consistent with the evaluation metrics obtained during model validation, demonstrating the model's robustness and efficiency.

## 9. Conclusion

The invoice data extraction project successfully demonstrates the use of machine learning to handle diverse invoice formats and extract key information without hardcoded labels. The optimized model is efficient and performs well on client desktops, making it a practical solution for real-world applications.