Automated Label-Value Extraction from Multilingual Documents:

A Deep Learning Approach Using Optical Character Recognition

elm19-others

August 10, 2025

Abstract

This report presents the first phase of developing an intelligent system for automated extraction of label-value pairs from multilingual documents including forms, invoices, and various administrative documents. Our approach leverages advanced Optical Character Recognition (OCR) technologies to process documents in French, Arabic, and English languages. The system aims to extract textual content with spatial positioning information and subsequently store structured data in a database. This phase focuses on comparative analysis of OCR engines, specifically Tesseract and PaddleOCR, evaluating their performance across different languages and document complexities.

1 Introduction

The digitization of paper-based documents and the extraction of structured information from unstructured formats represents a critical challenge in modern information management systems. Traditional manual data entry processes are time-consuming, error-prone, and economically inefficient. This project addresses the need for an automated solution capable of processing multilingual documents containing various forms, invoices, and administrative papers.

Our objective is to develop an AI-powered system that can:

- Extract textual content from images and PDF documents
- Handle multilingual content (French, Arabic, English)
- Identify spatial relationships between labels and their corresponding values
- Store structured data efficiently in a database

This report documents Phase 1 of the project, focusing on text extraction methodologies and OCR engine evaluation.

2 Methodology

2.1 OCR Engine Selection and Evaluation

Given the multilingual requirements of our project, we evaluated two prominent OCR solutions: Tesseract OCR and PaddleOCR. Both engines provide comprehensive language support and are well-documented with active community support.

2.1.1 Tesseract OCR

Tesseract, developed by Google, is an open-source OCR engine that supports over 100 languages. It utilizes Long Short-Term Memory (LSTM) neural networks for character recognition.

Listing 1: Tesseract Implementation Example

```
import pytesseract
   from PIL import Image
2
   import cv2
3
4
   def extract_text_tesseract(image_path, language='fra+ara+eng'):
5
6
       Extract text using Tesseract OCR
7
8
       # Load and preprocess image
9
       image = cv2.imread(image_path)
10
       gray = cv2.cvtColor(image, cv2.COLOR_BGR2GRAY)
11
12
       # Apply OCR
13
       extracted_text = pytesseract.image_to_string(
14
15
            gray,
            lang=language,
16
            config='--psm<sub>□</sub>6'
17
18
19
       return extracted_text
20
21
22
   # Example usage
23
   text = extract_text_tesseract('sample_document.jpg')
   print(text)
24
```

2.1.2 PaddleOCR

PaddleOCR is a practical ultra-lightweight OCR system developed by PaddlePaddle. It supports detection and recognition of natural scenes and document texts in multiple languages with high accuracy.

Listing 2: PaddleOCR Implementation Example

```
from paddleocr import PaddleOCR
   import cv2
2
   import numpy as np
3
4
   def extract_text_paddle(image_path, language='multilingual'):
5
6
       Extract text using PaddleOCR with bounding box information
7
8
       # Initialize OCR
       ocr = PaddleOCR(use_angle_cls=True, lang=language)
10
11
^{12}
       # Perform OCR
       result = ocr.ocr(image_path, cls=True)
13
14
       extracted_data = []
15
       for line in result:
16
            for word_info in line:
17
                bbox = word_info[0]
                                     # Bounding box coordinates
18
                text = word_info[1][0] # Extracted text
19
                confidence = word_info[1][1] # Confidence score
20
^{21}
                extracted_data.append({
22
                    'text': text,
23
```

```
'bbox': bbox,
24
                     'confidence': confidence
25
                })
26
27
       return extracted_data
28
29
30
   # Example usage
   data = extract_text_paddle('sample_document.jpg')
31
   for item in data:
32
       print(f"Text:_{\text'}]},_\Confidence:_{\text'}\; confidence']:.2f}")
33
```

2.2 Comparative Analysis Results

An analysis of online reviews and technical specifications for both OCR engines reveals some notable distinctions:

| Criteria | Tesseract | PaddleOCR |
|-----------------------|-----------------------|-----------|
| Printed Text Accuracy | High | Very High |
| Handwritten Text | Limited | Better |
| Complex Layouts | Moderate | Excellent |
| Multilingual Support | Good | Excellent |
| Bounding Box Info | Limited | Native |
| Processing Speed | Fast | Moderate |

Table 1: OCR Engine Comparison

3 Implementation and Results

3.1 PaddleOCR: The Chosen Solution

Based on our comparative analysis, PaddleOCR emerged as the superior choice for our requirements. The decision was primarily driven by:

- 1. **Superior handling of complex layouts**: Documents such as invoices and forms often contain tables, multiple columns, and varied text orientations.
- 2. **Better handwritten text recognition**: Many forms contain handwritten elements that Tesseract struggles to process accurately.
- 3. Native bounding box extraction: Critical for our subsequent label-value pairing algorithms
- 4. Robust multilingual capabilities: Seamless processing of mixed-language documents.

3.2 Bounding Box Extraction and Spatial Analysis

PaddleOCR's architecture follows a two-stage approach:

- 1. Text Detection: Identifies text regions in the image using a detection model
- 2. **Text Recognition**: Recognizes characters within detected regions using a recognition model

The bounding box coordinates provided by PaddleOCR are essential for our project as they enable spatial relationship analysis between labels and values without requiring additional computer vision models.

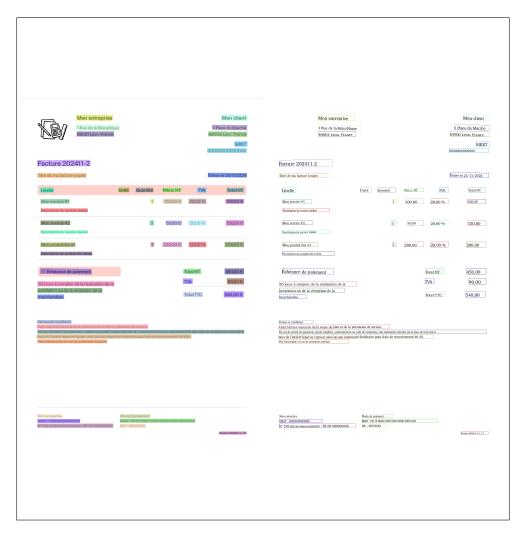


Figure 1: Sample Document with PaddleOCR Bounding Box Detection(more clear version in our code

3.3 Language-Specific Processing

3.3.1 French and English Processing

For Latin-script languages (French and English), PaddleOCR performs exceptionally well with minimal preprocessing required.

3.3.2 Arabic Text Processing

Arabic text processing presented unique challenges due to its right-to-left (RTL) reading direction. While PaddleOCR correctly identifies Arabic characters, the default output maintains left-to-right ordering. We implemented a post-processing function to handle this:

Listing 3: Arabic Text Direction Correction

```
import arabic_reshaper
from bidi.algorithm import get_display

def process_arabic_text(text):
```

```
"""Process Arabic text for proper display (right-to-left)"""
5
6
       trv:
         parts = text.split()
7
8
         fixed_parts = []
9
         for part in parts:
10
11
              # Reverse Arabic-like parts only (heuristic: check if contains Arabic
                 letters)
              if any('\u0600' <= c <= '\u06FF' for c in part):
12
                  part = part[::-1] # Reverse characters in the part
13
             fixed_parts.append(part)
14
15
         fixed_text = "".join(fixed_parts)
16
         # Reshape Arabic text for proper display
17
         reshaped_text = arabic_reshaper.reshape(fixed_text)
18
19
         # Apply bidirectional algorithm
20
         display_text = get_display(reshaped_text)
^{21}
         return display_text
22
23
           # If reshaping fails, return original text
           print("failed")
24
           return text
25
```

4 Current Challenges and Limitations

4.1 Multi-language Document Processing

Currently, there is no single OCR configuration that optimally handles all three target languages simultaneously. For documents containing mixed languages or when the primary language is unknown, we propose implementing a confidence-based selection mechanism:

- 1. Execute OCR with each language configuration (French, Arabic, English)
- 2. Compare confidence scores for extracted text
- 3. Select results from the configuration with highest overall confidence
- 4. Merge results where appropriate

4.2 Domain-Specific Text Recognition

The nature of our target documents may include:

- Domain-specific terminology
- Handwritten annotations
- Low-quality scanned images
- Irregular formatting

These factors may necessitate fine-tuning PaddleOCR models on domain-specific datasets. However, this approach presents challenges:

- 1. Data Collection: Acquiring sufficient annotated training data for specialized domains
- 2. Model Complexity: PaddleOCR's architecture requires careful consideration for fine-tuning
- 3. Arabic Language Optimization: Limited resources for Arabic OCR fine-tuning compared to Latin scripts

5 Future Work

Phase 2 of the project will focus on:

- 1. Label-Value Pair Extraction: Developing algorithms to identify relationships between text elements using spatial positioning
- 2. Multi-language Confidence Integration: Implementing the proposed confidence-based language selection system
- 3. Database Integration: Designing and implementing structured data storage solutions
- 4. Model Fine-tuning: If required, fine-tuning OCR models on domain-specific datasets
- 5. **Performance Optimization**: Enhancing processing speed and accuracy for production deployment

6 Conclusion

Phase 1 of our project successfully established a robust foundation for multilingual document text extraction. PaddleOCR emerged as the optimal solution, providing superior accuracy for complex document layouts while offering essential spatial information through bounding box coordinates. The system currently handles French, Arabic, and English texts with appropriate post-processing for RTL languages.

The identification of current limitations, particularly in multi-language document processing and domain-specific optimization, provides clear direction for subsequent development phases. The spatial information captured during text extraction positions the project well for the next phase of automated label-value pair identification.

Acknowledgments

We acknowledge the open-source communities behind Tesseract OCR and PaddleOCR for their invaluable contributions to optical character recognition technology.