

OCR Data Extraction Report

1. Introduction

This report details the functionality and techniques employed in the provided Python script for Optical Character Recognition (OCR) based data extraction from ID card images. The script is designed to automate the process of reading key information such as Name, Father Name, Gender, Identity Number, Dates (Birth, Issue, Expiry), and Country of Stay, even from images captured under varying conditions (e.g., distant or close-up).

2. Key Techniques and Their Working

The script leverages several advanced image processing and OCR techniques to achieve robust data extraction:

2.1. Image Normalization and Perspective Correction

- **Technique:** Gaussian Blur, Canny Edge Detection, Contour Finding, Perspective Transformation.
- **Working:** The script first preprocesses the input image by applying a Gaussian blur to reduce noise and then uses Canny edge detection to highlight the boundaries of objects. It identifies the dominant four-sided contour, which is assumed to be the ID card. If a card is detected from a distance or at an angle, a perspective transformation is applied. This "straightens" the card and scales it to a uniform size (e.g., 1600x1000 pixels). This standardization is critical as it ensures that regardless of the original photo's angle or distance, the card's appearance is consistent for subsequent OCR steps. If no clear card contour is found (implying a close-up image), the entire image is simply resized to the target dimensions.

2.2. Image Preprocessing Pipeline

- **Techniques:** Pixel Intensity Normalization, Non-Local Means Denoising, Morphological Operations (Opening and Closing).
- **Working:** After normalization, the image undergoes further enhancement:
 - **Pixel Intensity Normalization:** Adjusts the brightness and contrast of the image to ensure pixel values span the full dynamic range, making text more discernible.
 - **Noise Removal:** Employs a fastNlMeansDenoising filter to intelligently remove image noise without blurring important text details.
 - **Thinning/Skeletonization:** Uses morphological operations to refine text

strokes, aiming for more uniform character thickness. This can improve Tesseract's ability to recognize characters.

2.3. Region of Interest (ROI) Based Extraction

- **Technique:** Predefined Coordinate-based Cropping.
- **Working:** Instead of attempting OCR on the entire image, the script relies on a CARD_ROIS dictionary. This dictionary holds precise (x, y, width, height) coordinates for each specific data field (e.g., Name, Identity Number) on the *normalized* ID card. The script crops these exact regions from the preprocessed image and performs OCR only on these smaller, targeted areas. This significantly improves accuracy and efficiency by focusing the OCR engine on relevant text.

2.4. Multi-Configuration OCR with Tesseract

- **Technique:** pytesseract with varied Page Segmentation Modes (--psm) and OCR Engine Modes (--oem).
- **Working:** The script uses pytesseract to interface with the powerful Tesseract OCR engine. It employs a strategy of trying multiple Tesseract configurations. Different --psm values (e.g., psm 3 for general text blocks, psm 8 for single words, psm 10 for single characters) are used to adapt to the varying layouts of different fields. Crucially, it prioritizes the modern Neural Net LSTM engine (--oem 1) for better accuracy. By testing various combinations, the script selects the result with the highest confidence and longest text, ensuring the best possible OCR outcome for each ROI.

2.5. Post-OCR Text Cleaning

- **Technique:** Regular Expressions (re module).
- **Working:** Raw OCR output often contains noise, extra characters, or misinterpretations. The `clean_text` function applies field-specific regular expressions to refine the extracted text. For example:
 - For "Identity Number", it removes any leading non-digit characters and ensures only digits and hyphens remain.
 - For "Country of Stay", it strips common prefixes like "Country of Stay" or "SQuntry of Stay" and ensures the result starts with an alphabet.
 - Dates are cleaned to contain only numbers and common date separators.
 - Gender is standardized to 'Male' or 'Female' from various possible OCR outputs.

3. Script Workflow

The data extraction process follows these sequential steps:

1. **Initialization:** The script starts by loading the target ID card image.
2. **Image Preprocessing:** The OCRPreprocessor normalizes the image (correcting perspective and scaling) and applies a series of enhancements (pixel normalization, noise removal, morphological operations) to prepare it for optimal OCR.
3. **ROI-Based Text Extraction:** The ROIExtractor takes the preprocessed image and, for each defined Region of Interest (ROI), crops that specific area.
4. **OCR Application:** The cropped ROI is passed to the OCRTextExtractor, which runs multiple Tesseract configurations to extract the most confident text.
5. **Text Cleaning:** The extracted text for each field undergoes a specific cleaning process to remove noise and standardize formats.
6. **Data Consolidation:** The CardDataExtractor compiles all the cleaned, extracted fields into a structured dictionary.
7. **Results Display & Storage:** The final extracted data is printed to the console, showing a success rate, and is also saved to an Excel file (idcard_data.xlsx) for record-keeping.

4. How to Use the Script

To utilize this OCR script:

1. **Install Dependencies:** Ensure all required Python libraries are installed: opencv-python, pytesseract, numpy, Pillow, and openpyxl. You can install them via pip:

```
pip install opencv-python pytesseract numpy Pillow openpyxl
```
2. **Install Tesseract OCR Engine:** This script relies on the external Tesseract OCR engine. Download and install the Tesseract executable for your operating system (e.g., from [Tesseract at UB Mannheim](#)). Ensure its installation path is added to your system's PATH environment variable, or manually specify it in your script if needed.
3. **Update IMAGE_PATH:** In the main() function of the script, modify the IMAGE_PATH variable to point to the actual location of your ID card image file.
4. **Tune CARD_ROIS (Crucial for Accuracy):** The CARD_ROIS dictionary in the main() function defines the precise rectangular areas for each data field on the *normalized* ID card. You will need to determine these coordinates accurately for your specific ID card layout.
 - Initially, you can set SHOW_PREPROCESSING_STEPS = True to visually inspect the "Original Image with Detected Contour" (if applicable) and the "Final Preprocessed Image".

- Use an image editing tool to get the (x, y, width, height) pixel coordinates for each field on these *normalized* images.
- Update the CARD_ROIS dictionary with these values. This step is iterative and vital for achieving high accuracy.
- Once tuned, set SHOW_PREPROCESSING_STEPS = False to run the script without image pop-ups.

5. **Execute the Script:** Run the Python file from your terminal.

5. Demo Report: Extracted Data Example

Below is an example of the structured data that the script aims to extract and display. The actual values will depend on the input image and the accuracy of the OCR process.



FINAL EXTRACTED CARD DATA

=====

- ✓ Name :
- ✓ Father Name :
- ✓ Gender :
- ✓ Country of Stay:
- ✓ Identity Number:
- ✓ Date of Birth :
- ✓ Date of Issue :
- ✓ Date of Expiry :

=====



Success Rate: 8/8 fields (100.0%)

- ✓ Data saved to: idcard_data.xlsx

6. Drawbacks and Considerations

While this OCR script employs robust techniques, its performance can be significantly impacted by the quality and characteristics of the input images. Key drawbacks and considerations include:

- **Background Contrast:** The script's ability to accurately detect the ID card (especially for perspective correction) is **highly dependent on the contrast between the ID card and its background**. If the ID card is placed against a **lighter background**, particularly a white or very light-colored surface, it can be challenging for the contour detection algorithm to reliably distinguish the card's boundaries. This can lead to inaccurate normalization or the system defaulting to

processing the entire image, potentially reducing OCR accuracy.

- **Image Capture Angle:** For optimal performance, it is **crucial that the image of the ID card is taken from directly above**. Capturing the image from the sides, even with the perspective correction implemented, can introduce distortions that are difficult to fully compensate for. Extreme angles can lead to:
 - **Loss of Detail:** Text might become too skewed or compressed, making it unreadable even after correction.
 - **Inaccurate ROI Mapping:** While the script normalizes the card, severe initial distortions can still affect how precisely the predefined ROIs align with the actual text areas on the corrected image.
 - **Shadows and Glare:** Side angles are more prone to casting shadows or creating glare across the card, obscuring text and drastically reducing OCR accuracy.

7. Conclusion

This OCR script provides a robust and adaptable solution for extracting structured data from ID card images. By combining advanced image preprocessing techniques, a targeted Region of Interest (ROI) approach, multi-configuration OCR with the powerful Tesseract engine, and intelligent post-OCR text cleaning, it aims to achieve high accuracy even with varying input image qualities. The modular design, with distinct classes for preprocessing, ROI extraction, OCR, and data handling, makes the code maintainable and extensible. While the CARD_ROIS require initial tuning for specific ID card layouts, this investment ensures precise and reliable data extraction, making the script a valuable tool for automated document processing. Understanding the limitations related to background contrast and image capture angle is important for maximizing the script's effectiveness in real-world applications.