# Multi-Stage Framework to Boost Optical Character Recognition Performance on Low Quality Document Images

Nitin Gupta, Shashank Mujumdar, Abhinav Jain, Douglas Burdick

IBM Research

Document Intelligence Workshop at KDD 2021

# Improving OCR on Low-Quality Document Images - Overview

Factors affecting **Tesseract Engine**'s ability to generate binary representations and perform page segmentation :-
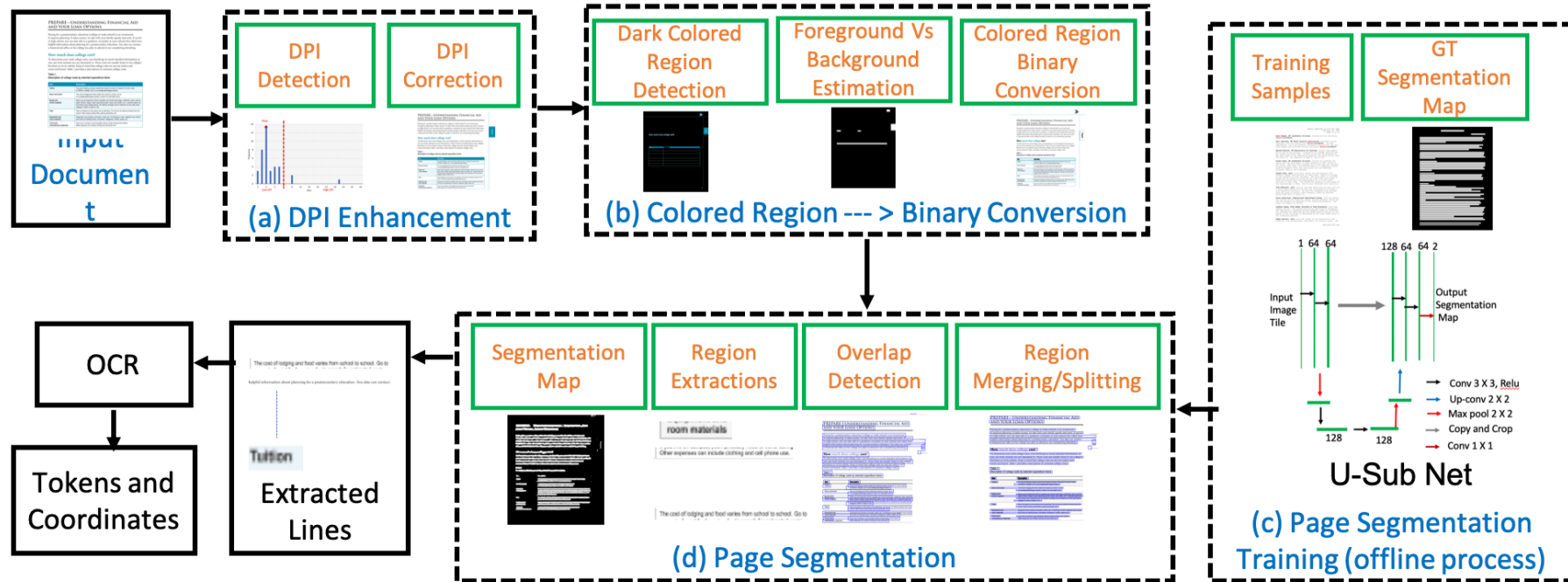
- Low Resolution

- Illumination Change

- Blur

- Noise

- Character Merging or Fragmentation

- Colored Regions with poor text contrast against background


We argue that any approach catering a specific issue in low-quality documents is sub-optimal in improving OCR performance. Thus, we propose a multi-stage framework that independently improves the performance of Tesseract at every stage.

Lastly, we present results on five challenging document image datasets and show superior performance against state-of-the-art baselines.

Document Intelligence Workshop at KDD 2021

# Proposed Pipeline Overview



(a) DPI Enhancement

- Input Document
- DPI Detection
- DPI Correction

(b) Colored Region ---> Binary Conversion

- Dark Colored Region Detection
- Foreground Vs Background Estimation
- Colored Region Binary Conversion

(c) Page Segmentation Training (offline process)

- Training Samples
- GT Segmentation Map
- U-Sub Net
  - Input Image Tile
  - Output Segmentation Map
  - 1 64 64 128 64 64 2
  - 128 128
  - Conv 3 X 3, Relu
  - Up-conv 2 X 2
  - Max pool 2 X 2
  - Copy and Crop
  - Conv 1 X 1

(d) Page Segmentation

- Segmentation Map
- Region Extractions
- Overlap Detection
- Region Merging/Splitting

- OCR
- Tokens and Coordinates
- Extracted Lines

Our framework consists of:–

- **DPI Enhancement Module** to identify and up-scale low resolution document images
- **Colored Region Detection** Module to detect and binarize colored text regions
- **Page-Segmentation Module** to extract text lines

Document Intelligence Workshop at KDD 2021

# DPI Enhancement

- Tesseract works best on document images with >300 DPI but suffers on 72-100 DPI low quality scans.

- We intend to re-scale the low 72-100 DPI images to 300 DPI using bi-cubic interpolation.

- From our comprehensive analysis, we identify that text lines in low 72-100 DPI images are of pixel height no more than 20 pixels. Thus, using this as a threshold, we detect and interpolate 72-100 DPI images.

# Colored Region – Binary Conversion

- Existing OCR systems use Otsu Thresholding on the whole image to achieve binarization which fails on regions containing text against colored background such as headers, headings highlighted text etc.

- The local contrast observed between text and colored-regions is different than what is observed in rest of the document with mostly dark-text appearing against light background.

- We contain those local contrasts by detecting colored regions as clusters (use K-Means) in the HSV space rep[...] the document and locally perform Otsu-thresholding in each detected region to achieve binarization.

Using Elbow method → **Find K$_{optimal}$**
Find colors in the image (HSV space) → **Clustering using K$_{optimal}$**
Find HSV range to represent each color → **Find Quantile 1 and Quantile 3**
→ **Extract every color patch**

Colored Region Extraction Module

(a) Original Image  (b) Binary  (c) OTSU  (d) Adaptive Mean  (e) Adaptive Gaussian  (f) Our Color Region Detection

No Text Detected from Header — No Text Detected from Header — No Text Detected from Header — Following two tokens detected: [Cia; tertyeatiny|- — All tokens detected correctly: Country; New; Yearly; Salary; average; Country; New; Yearly; Salary; average;

# Page Segmentation

- With complexity and variations of 2D layouts, page segmentation is challenging

- We formulate the page segmentation as image-to-image transformation problem.

- We motivate the required architecture from U-Net, used in medical image segmentation. The output is a per-pixel probability estimate that the pixel is a part of segmented text lines. We set the threshold of the output to get the binary segmented image.

- We further perform post-processing like Region smoothing, Region Overlap Detection and Region Merg... tighter bounding boxes in the segmented images.

Visual Illustration of Tesseract (a and c) and Proposed (b and d) Page Segmentation Pipeline on 72 D

(a)  (b)  (c)  (d)

# Results

Compare the OCR accuracy by matching the tokens extracted from the processed image with those from GT at the corresponding location returned by Tesseract Engine

|  | Tesseract | SRCNN | SAE-18 | Model_P | Model_PH | Model_PD | Model_PHD |
|---|---|---|---|---|---|---|---|
| LCWA | **95.12** | 91.95 | 92.78 | 94.99 | 94.41 | 94.92 | 94.54 |
| ICDAR | **95.09** | 92.27 | 91.13 | 94.68 | 95.08 | 94.66 | 94.72 |
| UNLV-A | 87.29 | 79.41 | 84.65 | **88.16** | 87.56 | 87.78 | 88.10 |
| UNLV-B | 94.40 | 80.77 | 93.41 | **94.54** | 94.28 | 94.13 | 94.47 |
| CI | 93.32 | 91.27 | 90.42 | 93.60 | 94.11 | 94.49 | **95.10** |

Table 1: OCR Accuracy Percentage on 300 DPI Images.

|  | Tesseract | SRCNN | SAE-18 | Model_P | Model_PH | Model_PD | Model_PHD |
|---|---|---|---|---|---|---|---|
| LCWA | 63.17 | 7.5 | 9.52 | 81.16 | 85.67 | 87.79 | **87.89** |
| ICDAR | 59.07 | 4.51 | 5.67 | 78.16 | 78.22 | 84.89 | **85.13** |
| UNLV-A | 18.76 | 6.43 | 10.28 | 32.46 | 32.74 | 33.70 | **33.76** |
| UNLV-B | 35.27 | 3.8 | 4.6 | 42.99 | 44.0 | 45.01 | **45.16** |
| CI | 57.63 | 12.92 | 29.71 | 73.95 | 74.71 | 80.43 | **86.60** |

Table 2: OCR Accuracy Percentage on 72 DPI Images.

Document Intelligence Workshop at KDD 2021

# Thank You