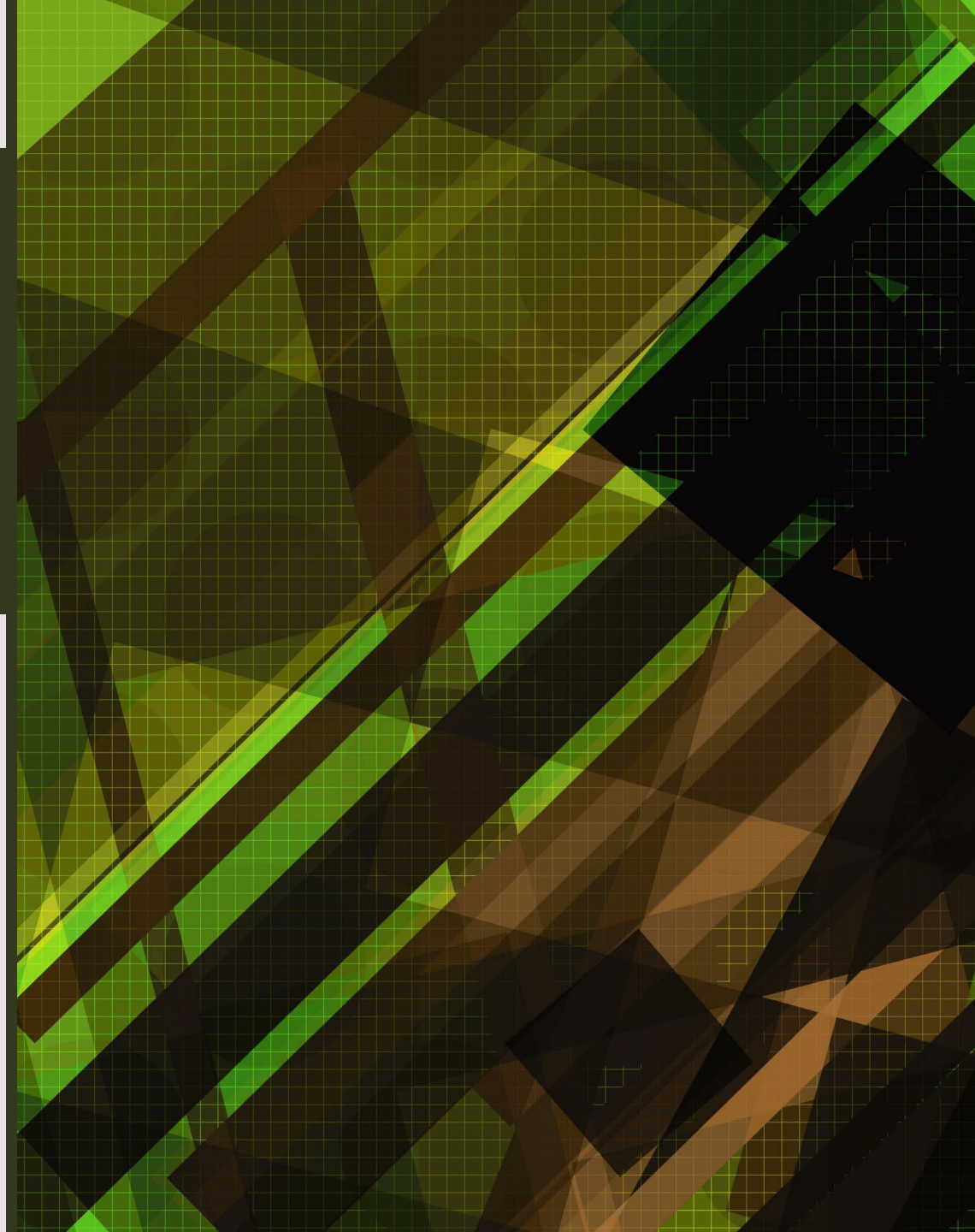# USING DEEP LEARNING TECHINQUES TO IMPORVE OPTICAL CHARACTER RECOGNITION FOR HISTORICAL DOCUMENTS

by Mengyue Zhao

beczhaozmy@gmail.com

# Problem Introduction

Historians have a hard time extracting large amounts of data from raw scans. Especially from historical tables which commercial for which commercial software are not designed.

Linking different historical sources requires very high accuracy of picture-to-data-frame recognition.

I found a very interesting large dataset that requires an improvement on such software. This project is designed to handle tens of thousands of raw scans as shown in the picture on the right side.

The alternative to this solution is manual input by typists, which is costly, inaccurate and often-time yields unknown or stringent mistakes that are hard to correct.

Tesseract is an optical character recognition software initially maintained by HP in the 1990s, later open-sourced and now maintained by Google.

# What is Tesseract and how does it work?

# How does Tesseract (or OCR in general) Work?



Pre-processing

Page Layout Analysis

Text Line Analysis

Character Segmentation

Optical Character Recognition (OCR)

Post Processing

# Tesseract Architecture

**Page Layout Analysis***

Output: Text Regions/ Block of Texts

→

**Connected Component Analysis**

Output: Text Lines

**Chopping with Character Spacing**

Output: Segmented Words

→

**First Pass with a CNN or LSTM**

Output: recognized characters, high confidence words are used as training data for the second pass

**Second Pass with an Adaptive Classifier Using LSTM or CNN**

Output: proposal of prototype words

→

**Linguistic Analysis**

# Page Layout Analysis

## Bottom-up Approach

| Connected Component Analysis on pixels | Morphological transformation | Combine regions where possible |



Connected Component Analysis

Text Components | Candidate Tab-Stop Components | Detected Tab Segments

## Top-down Approach

| Proposed Page Layout (pre-trained model) | Cut Images into the proposed parts | Validate with Confidence Level |



Block Finding

Detected Columns | Typed Column Partitions | Detected Blocks

# Using LSTM in OCR



Figure 1. The network architecture. The architecture consists of three parts: 1) convolutional layers, which extract a feature sequence from the input image; 2) recurrent layers, which predict a label distribution for each frame; 3) transcription layer, which translates the per-frame predictions into the final label sequence.

# Links

**Visualized Tesseract Architecture**

https://limitlessdatascience.wordpress.com/2019/07/01/tesseract-4-0-intro-installation/

**How does Tesseract use LSTM**

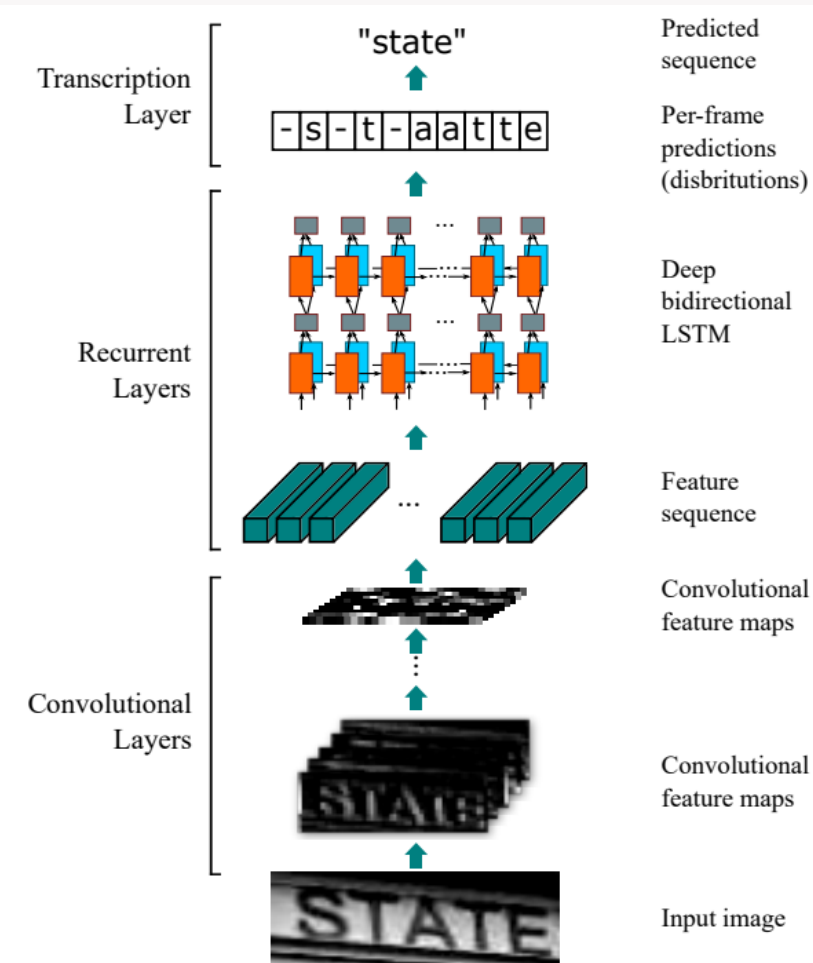https://github.com/tesseract-ocr/docs/blob/master/das_tutorial2016/6ModernizationEfforts.pdf¥

https://nanonets.com/blog/ocr-with-tesseract/

**How LSTM work for OCR**

https://towardsdatascience.com/a-gentle-introduction-to-ocr-ee1469a201aa

# Page Segmentation Solutions

- Left: Tesseract Segmentation Results

- Right: CV Morphological Transformations

# OCR Results

# **Proposed Solutions**

- Customized Preprocessing [done]

- Page Layout Analysis with Morphological Transformation [done]

- Page Segmentation with Morphological Transformation [done]

- Fixed-width Character Segmentation [tuning]

- Retrain Neural Network on Specific Fonts [done]

- Restrained Classifier with a Customized Dictionary (word list) [doing]

- Restricted Output Base [done]

- Post Processing [doing]