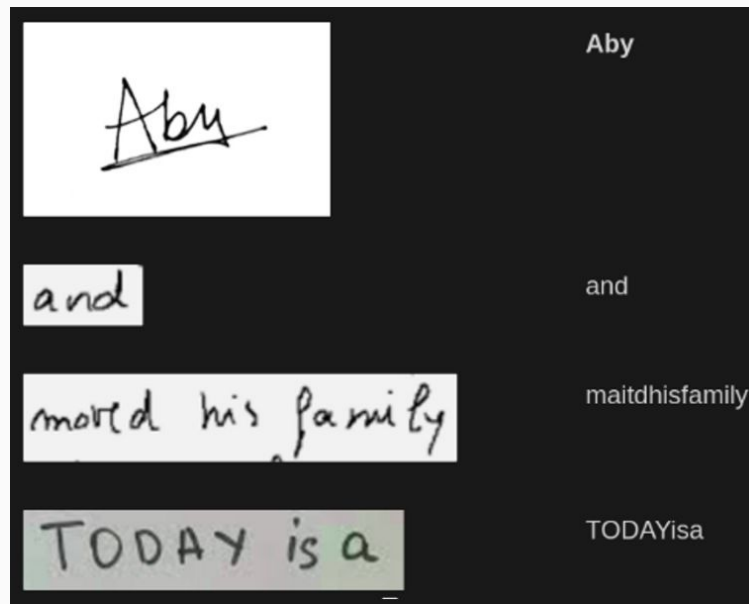


Improving OCR Pipeline @Decimal Point Analytics



Issue at hand.

- Current pipeline not end-to-end
 - Detection + Recognition
- Issue with recognition pipeline??
 - **Unable to reproduce white-spaces**
 - Unable to recognize special characters
 - Unable to parse numbers in financial amounts and hyperlinks
- Another model text-fuse-net is used to detect white spaces**



Why?

t permitted for pu

220 Manchester B

2371 Jerrold Ave

2371 Jerrold Ave

State Route 27

394 Manchester B

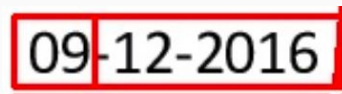
Approach #1: Improve detection to one word

- Baseline various word detector models:
 - WordNN, DBNet, CRAFT, Textfusenet, EasyOCR, Donut, LMV3, TrOCR
 - Markers: inference time and accuracy
 - Recurrent issues?
 - Missing out of small words
 - Merging of 2 or more words
- Dataset reasoning
 - Impact of training data
 - Custom training ??

Apart for this, I reviewed literature on text detection and followed current works exhaustively.



n B Jr



09-12-2016



Ca darera.



Post Box Services Plus

Approach #1: Improve detection to one word

match based (euclidean distance metric).
Brute force matching can be computationally expensive.
are the features in the image match.

- **K-nearest neighbor (KNN) matching** (Figure 4.10) is a fast and efficient algorithm for matching features. KNN matching is a two-step process: first the descriptors find the K nearest neighbors of each feature descriptor. Then, if there are the matches that are not unique. For example, if the nearest neighbor of a feature descriptor in one image is also the nearest neighbor of another feature descriptor in the other image, then the matches considered unique and kept (rest are discarded).

- **FLANN (Fast Library for Approximate Nearest Neighbors) matching** is a fast and efficient algorithm for matching features in large datasets. FLANN is a library that implements a collection of algorithms optimized for fast nearest neighbor search in large datasets and for high-dimensional features. It contains a collection of algorithms for searching nearest neighbor queries (KNN algorithm). These algorithms can be used to search for nearest neighbors within query point sets of points with arbitrary sizes. The algorithm in FLANN can also be used to find a good approximation of the nearest neighbors with a query point sets of points with arbitrary size.

2.3 Geometric Transformation

We can find the matches key points to estimate the homography matrix that aligns the images. However, why columns that in a given image stitching are determined the transformation (required) warp the images into common geometry (scale). Moreover, in real world scenarios, such as photography from drone images, the image of viewpoint and perspective is rarely, which make this step extremely necessary. Given the matched feature points, the geometric transformation between the two images is estimated using a robust estimation algorithm, such as Ransac (Random Sample Consensus) (RANSAC). The transformation can be a homography or affine transformation depending on the number of matching points and their distribution.

2.4 Image Blending



Figure 3: Image Blending with warp perspective

Once the homography matrix is estimated, we can use it to warp one of the images onto the plane of the other image, so that the two images are aligned. However, the resulting image may contain visible seams or artifacts due to differences in brightness, color, or texture between the two images.

We can use image blending techniques to smooth out the seams and create a more seamless transition between the images to address this issue. One common approach is to use a weighted average of the pixel values in the overlapping region, where the weights depend on the distance of each pixel from the boundary of the overlapping region. Finally, image blending with warp perspective is utilized in the pipeline, which allows us to create seamless panoramas from multiple overlapping images. It requires accurate estimation of the homography matrix and careful selection of the blending technique to ensure a smooth and seamless transition between the images. On our end, we implemented alpha and gaussian blending techniques to achieve this.

1. **Alpha blending** Alpha blending is a process that takes the transparency of each pixel into account. The degree of transparency, or alpha value, is used to blend the colors of the foreground and background images or objects in a way that creates a smooth transition between them.

- i. Choose a masking value M (transparency level)
- ii. Final image: $M * \text{img1} + (1 - M) * \text{img2}$

2. **Gaussian blending** Gaussian blending, on the other hand, is a process of blending images or objects using a Gaussian distribution. In this method, a Gaussian filter is applied to each pixel, which smooths out the sharp edges and creates a softer transition between the images or objects being blended. The degree of smoothing is determined by the size of the filter kernel.

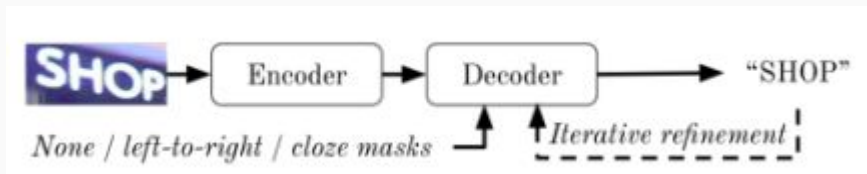
- i. Compute the Gaussian pyramid for each image
- ii. Compute the Laplacian pyramid for each image
- iii. Combine the left and right halves of each level of the Laplacian pyramid
- iv. Reconstruct the blended image from the Laplacian pyramid

3. **Seamless blending** This method uses an image processing operator that enables the merging of two images without creating any unsightly seams. This method also ensures that the color of the inserted image is likewise altered so that the inserted object appears to be a natural part of the target image's surroundings.

- Annotated 341 words
- Took 21 minutes
- Approximately 4 seconds / word
- Actual annotation time reported (1:26 was overhead)

Approach #2: Improve recognition phase

- What is PARSEQ?
 - Literature, set-up, training.
- Reviewed current state of text recognizers
- Impact of dataset?



- Issue at hand? How to create such a large custom dataset

“They are trained by enforcing an autoregressive (AR) constraint on the language context where future tokens are conditioned on past tokens but not the other way around”

Issues with dataset?

- Current checkpoint trained on **single word images**
- Lack of special characters
- Need to create a dataset which is rich in white-spaces and maybe other special characters; especially those from financial domain

Creating custom dataset

1. PyMuPDF;
 - a. Input : PDF
 - b. Output : Images of individual words; corresponding ground-truth text
 2. Tiger-Synth;
 - a. Input : Corpus of text
 - b. Output : Images with text as in input corpus
- Created **char_set**; bundled in **lmdb** format; set to **train**

Issues with dataset v0

- PyMuPDF: Random concat of words -> Model doesn't generalize
- Tiger-Synth: Invalid words -> Poisons language model
 - Can output ***par&eq, hell0, `good mourning`, `good dye`***

41.65%

Accuracy of model on test-set that was trained on v0 (136k training samples)

PS: Trained 3 more different models with minor tweaks but negligible trade-off gain

Dataset v1: `Semantic dataset`

- Words occur in order of natural language
- Individual words are valid english words

Outcome?

- Labels coming from spoken english gave good results
- White-spaces recognized
- Rest; unsatisfactory

Ignoring punctuation; accuracy:

88%

Hypothesis #1: Special characters under-representation

(PS: Hypothesis #2: overfitting was ruled out after experimentation)

(PS: Hypothesis #3: shallow model architecture was ruled out)

- 101 special characters
- Frequency of SC: **10.34%** (~ 100-88)
- 62.58% of labels had one or more SC
- 94/101 SCs are starved; their occurrence < 500; (out of total 10M characters)

Failure Categories:

- Confusing chars:
* vs *, ± vs +, `` vs ``
- Non-english Latin chars:
alpha, beta, gamma, ...

Next step?

Unicode normalization

```
In [1]: print("\u00C7", "\u0043\u0327")
        Ç Ç

In [2]: "\u00C7" == "\u0043\u0327"

Out[2]: False

In [3]: "Ç" == "Ç"

Out[3]: False
```

99.23%

Accuracy of model on test-set that was trained on v1.2 + ***partial unicode normalization***

PS: Tested on different population of dataset and save similar results

Refining edge cases

- Email IDs and hyperlinks
- Dates
- Function names
- HTML code and other programming languages
- Tags and mentions
- Mathematical equations
- IPs and amounts

Parallely, created a dataset with 1.4M labels (10x the original) using;

- **Gov docs**
- **NCERT books**
- **Research papers**
 - **Trade-off was however found unsatisfactory**

>Also explored the issue of superscripts and subscripts and how it cannot be tackled with the given setup.

> Inference, visualization, testing, demo, early-stopping, logging scripts

Deploying and documentation

- Installation
- Inferencing
 - With ground truth available
 - Without ground truth available
 - With image
 - Without image
- Training
- Creating custom dataset
- Model description and logs
- Dataset description

PARSEQ-robust

Getting started

This repository contains code for fine-tuning the OCR pipeline and its inference and other supporting scripts.

- To view unedited HTML-Logs (*recommended*); download the repo locally and open the html file [here](#)
- Alternatively, logs in PDF can be found [here](#)
- Please note; that logs run from bottom to top; that is, initial stages of development are down below and more recent developments were logged above.

Repo description

- [assets](#) : Supporting documents and scripts
- [ds-maker](#) : Script for constructing dataset
- [sub-model](#) : Model checkpoints and logs
- [suite-infer](#) : Demo script
- [v2](#) : Codebase for training/ testing

Image

€300.6 tax

meet.google.com/fwo-iyxe-xif

c.35,000

953,843 tCO₂e/89.1%

£30.3m,

(3) Third point

Group's effective tax rate has benefited in

\$90.3 dollars

okay-this-is-insanely-large-text

PARSEQ-old

300.6tax

meet.google.com/two-iy
xe-

35,000

953,843tCO₂e/891%

30.3m,

(3)Thirdpoint

Groupseffectivetatehac
din

\$90.3dollars

okay-this-insamely-lage
xt

PARSEQ-new

€300.6 tax

meet.google.com/fwo-iy
xe-xif

c35,000

953,843tCO₂e/89.1%

£30.3m,

(3) Third point

Group's effective
taxrate has benefited in

\$90.3 dollars

okay-this-is-insanely-lar
ge-text

Wrapping things up....

- Old model unable to detect white-spaces and certain characters
 - Hallucinates with long labels and symbol-heavy labels
- Time saved by using 1 model instead of two
- Additional cases, symbols, contexts, data-populations can easily be covered with the current codebase set-up.

Thank you
:)