

# Using Optical Character Recognition in an Applied Research Workflow

Rohan Alexander\*

John Tang<sup>†</sup>

Diego Mamanche Castellanos<sup>‡</sup>

14 December 2020

## Abstract

We discuss the use of Optical Character Recognition (OCR) as a initial step in an applied research workflow. OCR is a way of converting images, for instance in a PDF file, into text and numbers that can then be analysed. We review a variety of methods and their relative merit. The choice between these approaches turns on the scale of the OCR required, the level of funding that is available, and the user's technical skills. We provide a hueristic to guide this decision and worked examples. Finally, we place this within the context of an illustrative applied research workflow detailing its relationship with other aspects of the workflow.

## 1 Introduction

A crucial step in the data analysis workflow is gathering data. This step involves tools and methods to collect raw data in an systematic way. In this paper we are interested in gathering data from Portable Document Format (PDF) files, in particular, when these PDFs were created by scans, and hence must be processed before the data they contain can be prepared, cleaned, and analysed. Such PDF files are generally unstructured data, meaning that the data are not tagged nor organised in a regular way.<sup>1</sup> While tags in the form of a text replacement for each character image could be assigned manually, typically a statistical model is used for reasons of efficiency and reproducibility, and this process is called Optical Character Recognition (OCR). We evaluate several pre-trained commercial OCR options accessed through proprietary APIs: 1) Abbyy, 2) Amazon Textract, 3) Azure Cognitive Services, and 4) Google Vision; additionally several pre-trained and trainable OCR options that a user runs themselves: 5) Kindai OCR, 6) Kraken, and finally 7) Tesseract. We find... [TBD].

PDF files were developed in 1993 by the technology company Adobe and are useful for documents because they are meant to display in a consistent way independent of the environment that created them or the environment in which they are being viewed. A PDF viewed on an iPhone should look the same as on a Android phone, as on a Linux desktop. One feature of PDFs is that they can include a variety of objects, for instance, text, photos, figures, etc. However, this variety can limit the capacity of PDFs to be used directly as data inputs for statistical analysis. The data first needs to be extracted from the PDF.

Sometimes it is possible to copy and paste the data from the PDF. This is more likely when the PDF only contains simple text or regular tables. In particular, if the PDF has been created by an application such as

---

\*University of Toronto, rohan.alexander@utoronto.ca.

<sup>†</sup>University of Melbourne.

<sup>‡</sup>University of Toronto.

<sup>1</sup>Losee (2006) describes how structured data have a highly-regular organization or a there is some pre-defined data model where regularities apply to everything in a dataset, for instance a table. Losee (2006) further describes how semi-structured data contains this same information, but instead of a regular structure that applies to the entire dataset, there may be some structural information, for instance, tags such as ‘name = “Bob”’, or some other marker. Finally, unstructured data, such as texts or images, holds information with no explicit structured data. However, tags may be assigned using manual or automatic techniques, converting the unstructured data to semi-structured data. The diversity of possible data structures complicates the data gathering process as each may need a specific approach.

Table 1: OCR options

OCR Tool	Company	Type	Technical requirements
FineReader	Abbyy	Proprietary API	Low
Texttract	Amazon	Proprietary API	Medium
Azure Cognitive Services	Microsoft	Proprietary API	Medium
Vision	Google	Proprietary API	Medium
Kindai	n2i Project	Open source pre-trained OCR model	High
Kraken	Open source	Open source OCR engine	High
Tesseract	Open source	Open source OCR engine	Medium

Microsoft Word, or another document or form creation system, then often the text data can be extracted in this way because they are actually stored as text within the PDF. But it is not as easy if the text has been stored as an image which is then part of the PDF. This may be the case for PDFs produced through scans or photos of physical documents, and some older document preparation software. Furthermore, there is a linguistic component that adds more complexity (Kanagarathinam et al. 2019). For instance, non-Latin characters such as Japanese, Chinese, Arabic, among others, generate several difficulties in different stages of the text mining process. Not only do these languages have complicated structures and a large number of categories (for instance, compare the 26 characters in the English alphabet with the thousands of characters in Japanese kanji). But often many characters are similar and fonts may have a large effect. In those cases, the need for reliable approaches to replace an image containing many characters with text of those characters becomes important.

Optical character recognition (OCR) is a process that transforms a image of text into actual text. Although there may not be much difference to a human reading a PDF before and after OCR, the main difference is that the PDF is now machine-readable (Cheriet et al. 2007). OCR has been used to parse images of characters since the 1950s, initially using physical processes. Such approaches were understandably slow and have been replaced by a computer-based recognition system involving the following steps (Cheriet et al. 2007, 4):

- Pre-processing to both: increase the quality of the input image, and to focus on the data of interest.
- Feature extraction to identify the distinctive features of the characters to be recognized.
- The actual classification stage, in which those features are used to identify the characters.

In this paper we review various OCR options for data analysts who need to gather their data from PDFs where the characters of interest for analysis must first be transformed from image to text. The paper goes through the process of using various options to do this OCR and the tradeoffs that must be made. We consider a few different options in terms of the image that is to be transformed. We start with a one page extract of modern English, a one-page table of numbers, one page of modern Japanese, one page from an English book from the 1800s, and a page from a Japanese dictionary from the early 1900s.

We find... [TBD].

Our paper is similar to Rychlik et al. (2020) which focuses on Pashto, Farsi, and Traditional Chinese. Additional useful surveys include Lombardi and Marinai (2020) and Reul et al. (2019).

The remainder of this paper is structured as follows... [TBD]. Finally we make some suggestions for what is needed in the future.

## 2 Background

In this section we introduce the OCR options that we consider, their essential features, and cost. There are seven [UPDATE] OCR options evaluated in detail in this study (Table 1).

The gold-standard in terms of quality is a manual approach. For instance, a team of research assistants

who were paid to inspect each page, inputting the correct element and checking the inputs of others will likely provide the best recognition. While this option is only possible for smaller datasets and is not reproducible, sharing high-quality datasets created in this way is important for evaluating the output of automated approaches.

ABBYY - <https://www.abbyy.com> - is a private company that has conducted OCR through the product FineReader since 1993. FineReader has a variety of options that can be useful. There are a variety of options in terms of using FineReader including downloading an application locally, accessing it online, or using an iPhone app. The iPhone app can be especially useful when the item that is being scanned should not have its pages separated or even be fully flattened for scanning because it uses the camera. FineReader is capable of manual ‘training’ for instance, a user is able to identify unrecognised items. There are a variety of different prices depending on the application. FineReader online requires an annual subscription that allows for 5,000 pages for \$149 (USD), FineReader for Mac is \$150 (CAD).

Amazon Textract - <https://aws.amazon.com/textract/> - is an OCR application that is available through the Amazon. It can be accessed via an API or through Amazon Web Services (AWS). Textract is limited to only English for handwriting and even the printed text languages that Textract can detect are limited to English, Spanish, Italian, French, Portuguese and German (“Amazon Textract FAQs” 2020). Helpfully, the output can be put into a workflow such that low confidence recognitions can be automatically reviewed by a human. Textract requires a subscription and on a monthly basis the cost begins at \$1.50 per 1,000 pages for documents with text, and \$15.00 for documents with tables.

Microsoft Azure Cognitive Services provides OCR through the Computer Vision application - <https://azure.microsoft.com/en-us/services/cognitive-services/computer-vision/> (“Microsoft Azure Documentation” 2020). It is access via an API and handwriting OCR is only available for English, however printed OCR is available for a wide variety of languages including more complicated languages such as Arabic, Chinese, and Japanese. Mixed languages are also possible. The cost to use the API begins at \$1.50 per 1,000 pages.

Google Vision API - <https://cloud.google.com/vision> - is an OCR application for a wide variety of printed languages including those of Azure. It is available through an API. The cost to use the API begins at \$1.50 per 1,000 pages (“Google Vision API Documentation” 2020).

Kindai OCR (Le et al. 2019) - <https://github.com/ducanh841988/Kindai-OCR> - is a pre-trained model that uses an attention based encoder-decoder approach specifically designed for Japanese historical documents. It is part of a larger research program on Japanese character recognition including Horiuchi and Kato (2009) and Nguyen et al. (2017). This research program is complementary to research on historical Hanja (Korean) documents, for instance Kim et al. (2004), and historical Khmer (Cambodia) documents, for instance Valy, Verleysen, and Chhun (2020). Le, Clanuwat, and Kitamoto (2019) provides broader consideration of issues with the recognition of historical Japanese, and related, characters including ‘noised, damaged characters, and background’, as well as the possibility that characters may be ‘written in cursive and are connected’. The Center for Open Data in the Humanities - <http://codh.rois.ac.jp> - provides a large number of training datasets, and runs the n2i Project from which Kindai came out of.

OCRopus - <https://github.com/ocropus/ocropy> - is an open-source OCR engine initially developed in 2007 as ‘a Google-sponsored project’ (Breuel 2007). Kraken - <https://github.com/mittagessen/kraken> - is an ‘extensively rewritten fork of the OCROpus system’ by Benjamin Kiessling (Kiessling 2019). Martínek, Lenc, and Král (2020) use Kraken to create ‘an efficient domain-dependent OCR system that focuses on historical German documents by picking the best tools and approaches available’. Similarly, Romanov et al. (2017) use Kraken to achieve ‘accuracy rates for classical Arabic-script texts in the high nineties’.

Tesseract - <https://github.com/tesseract-ocr/tesseract> - is an open-source OCR engine initially developed ‘at HP between 1985 and 1995, shelved for 10 years, open-sources in 2006 and now developed mostly at Google’ (Smith 2013). It has had a variety of versions, with the latest being Tesseract 4 which added a ‘neural net (LSTM) based OCR engine’ (“Tesseract Manual Page” 2018). Tesseract is a widely used OCR engine with a variety of implementations available including the default usage of command line, Python through the pytesseract package [ADD CITE] and R through the Tesseract package (Ooms 2019) for which a useful starter example is Rodrigues (2019).

There are various other options that we considered, but ultimately did not evaluate in a serious way. The first of this is Amazon Rekognition - <https://aws.amazon.com/rekognition/> - which is a [ADD description]. However, the service only recognizes up to 50 characters per image. As this is too small for most OCR applications in academia we did not further explore this service.

## 3 Data

### 3.1 Book in English

We may be interested in the text from a book itself. In that circumstance, we'd often have scans of pages of each book, for instance as illustrated in Figure 1 illustrates a page containing a table from Brontë (1847).

Scanned pages of text from books will typically have some different fonts that need to be considered. They may also have page numbers and illustrations. Finally they may have various Pages of text Tables typically have several interesting features with regard to OCR. They typically have a page number and a title, both of which may or may not be important. They are then typically arranged in ordered rows and columns, which may or may not have lines separating them. Some of the cells may have symbols in them that may relate to footnotes at the end of the table. And finally, they may have additional comments and notes at the end of the table.

### 3.2 Table of numbers

A common reason for wanting to OCR a document is to use the tables that are contained in old books. Figure 2 illustrates a page containing a table from Bowley (1902)

Tables typically have several interesting features with regard to OCR. They typically have a page number and a title, both of which may or may not be important. They are then typically arranged in ordered rows and columns, which may or may not have lines separating them. Some of the cells may have symbols in them that may relate to footnotes at the end of the table. And finally, they may have additional comments and notes at the end of the table.

### 3.3 Japanese dictionaries from the early 1900s

[JOHN: Add a para about the origin of this dictionary.]

Here we focus on a page focused on 'a', which is in hiragana or in katakana Figure 3.

This page has several features of interest. Firstly, it is in vertical form and should be read from right to left. This means that the first five characters that we expect from an OCR output are: . The page is divided into four rows of content. An OCR workflow will need to take this into account during the pre-processing stage. In total, this page contains 3,293 characters, and while a majority of them are Japanese, there are still 99 English characters. Finally, there are also important nuances in terms of punctuation. For instance, certain types of brackets - [ADD AN EXAMPLE] - contain the origin of the word.

## 4 Results

### 4.1 Azure Cognitive Services - Computer Vision OCR

Among the Microsoft Azure's portfolio, Azure Cognitive Services offers Computer Vision, an API that processes images and returns information based on the visual features the user is interested in. Those services comprise object detection of an image, visual features tagging, image categorization, Optical Character Recognition (OCR), among *others*<sup>4</sup>. The OCR API processes an image and returns the language of the document, orientation, regions, lines, boundaries, and finally, the characters.

In Python we need several libraries for calling the API.

**I**T THERE was no possibility of taking a walk that day. We had been wandering, indeed, in the leafless shrubbery an hour in the morning; but since dinner (Mrs Reed, when there was no company, dined early) the cold winter wind had brought with it clouds so sombre, and a rain so penetrating, that further out-door exercise was now out of the question.

I was glad of it: I never liked long walks, especially on chilly afternoons: dreadful to me was the coming home in the raw twilight, with nipped fingers and toes, and a heart saddened by the chidings of Bessie, the nurse, and humbled by the consciousness of my physical inferiority to Eliza, John, and Georgiana Reed.

The said Eliza, John, and Georgiana were now clustered round their mama in the drawing-room: she lay reclined on a sofa by the fireside, and with her darlings about her (for the time neither quarrelling nor crying) looked perfectly happy. Me, she had dispensed from joining the group; saying, ‘She regretted to be under the necessity of keeping me at a distance; but that until she heard from Bessie, and could discover by her own observation that I was endeavouring in good earnest to acquire a more sociable and child-like disposition, a more attractive and sprightly manner—something lighter, franker, more natural as it were—she really must exclude me from privileges intended only for contented, happy, little children.’

‘What does Bessie say I have done?’ I asked.

‘Jane, I don’t like cavillers or questioners: besides, there is something truly forbidding in a child taking up her elders in that manner. Be seated somewhere; and until you can speak pleasantly, remain silent.’



Figure 1: A scanned page from a book.

## REVENUE OF THE UNITED KINGDOM.

Unit, in all columns, £10,000.

	Total Revenue.	Inland Revenue and Customs.	Customs.	Excise.	Property and Income Tax.	Post and Telegraph
1850	5,739	5,431	2,226	1,497	560*	216
1851	5,732	5,412	2,204	1,528	560*	228
1852	5,658	5,335	2,222	1,538	550*	237
1853	5,753	5,401	2,214	1,575	570*	237
1854	5,890	5,502	2,251	1,630	580*	252
1855	6,282	5,944	2,163	1,680*	1,070*	237
1856	7,026	6,601	2,324	1,730*	1,520*	281
1857	7,279	6,848	2,353	1,840*	1,620*	292
1858	6,788	6,309	2,311	1,782	1,159	292
1859	6,548	5,987	2,412	1,790	668	320
1860	7,109	6,570	2,446	2,036	960	331
1861	7,028	6,514	2,331	1,943	1,092	340
1862	6,986	6,412	2,367	1,833	1,036	351
1863	7,060	6,390	2,403	1,715	1,057	365
1864	7,021	6,306	2,323	1,821	908	381
1865	7,031	6,291	2,257	1,956	796	410
1866	6,781	6,036	2,128	1,979	639	425
1867	6,943	6,156	2,230	2,067	570	447
1868	6,960	6,204	2,265	2,016	618	463
1869	7,259	6,422	2,242	2,046	862	466
1870	7,543	6,708	2,153	2,176	1,004	477
1871	6,994	6,106	2,019	2,279	635	527
1872	7,471	6,484	2,033	2,333	908	543
1873	7,661	6,660	2,103	2,578	750	583
1874	7,734	6,608	2,034	2,717	569	700
1875	7,492	6,397	1,929	2,739	431	679
1876	7,713	6,525	2,002	2,763	411	719
1877	7,857	6,636	1,992	2,774	528	730
1878	7,774	6,610	1,997	2,746	582	746
1879	8,115	6,899	2,032	2,740	871	757
1880	7,934	6,695	1,933	2,530	923	777
1881	8,187	6,895	1,918	2,530	1,065	830
1882	8,396	7,058	1,929	2,724	994	863
1883	8,739	7,313	1,966	2,693	1,190	901
1884	8,616	7,187	1,970	2,695	1,072	947
1885	8,799	7,380	2,032	2,660	1,200	966
1886	8,958	7,493	1,983	2,546	1,516	989
1887	9,077	7,611	2,015	2,525	1,590	1,028
1888	8,980	7,566	1,963	2,562	1,444	1,060
1889	8,847	7,360	2,007	2,560	1,270	1,118
1890	8,930	7,341	2,042	2,416	1,277	1,177
1891	8,949	7,358	1,948	2,479	1,325	1,226
1892	9,099	7,534	1,974	2,561	1,381	1,263
1893	9,040	7,480	1,971	2,536	1,347	1,288
1894	9,113	7,543	1,971	2,520	1,520	1,301
1895	9,408	7,865	2,011	2,605	1,560	1,334
1896	10,197	8,512	2,076	2,680	1,610	1,422
1897	10,395	8,597	2,125	2,746	1,665	1,477
1898	10,661	8,855	2,180	2,830	1,725	1,518
1899	10,834	8,945	2,085	2,920	1,800	1,586

\* These figures cannot be given accurately within £100,000.

Figure 2: A scanned table from a book.



Figure 3: A page from an early 1900s Japanese dictionary.

```

import os
import sys
import requests
# If you are using a Jupyter notebook, uncomment the following line.
#%matplotlib inline
import matplotlib.pyplot as plt
from matplotlib.patches import Rectangle
from PIL import Image
from io import BytesIO
import matplotlib
from os import path
import ast

```

Then, we need to create headers and parameters necessary for calling the API. The library `requests` is commonly used in python to do that. The method `post` is needed as the API uses the HTTP method POST.

```

#Create the URL for the API call. The endpoint corresponds to the Azure's
# endpoint of the account.
ocr_url = endpoint + "vision/v3.0/ocr"

#Create the header using the Azure's subscription key.
headers = {'Ocp-Apim-Subscription-Key': subscription_key}

#In params language': 'ja' for japanese only. 'unk' means self detection
params = {'language': 'unk', 'detectOrientation': 'true'}
data = {'url': image_url}

#Call and save the response using the method post from the library request
response = requests.post(ocr_url, headers=headers, params=params, json=data)

```

Here we have an example of the response:

```

{'language': 'ja',
'textAngle': 0.0,
'orientation': 'Up',
'regions': [{}'boundingBox': '31,30,1521,2721',
'lines': [{}'boundingBox': '730,30,44,2718',
'words': [{}'boundingBox': '735,30,34,34', 'text': ' '},
{'boundingBox': '736,71,34,26', 'text': ' '},
{'boundingBox': '749,100,6,11', 'text': ' '},
{'boundingBox': '735,116,34,27', 'text': ' '},
{'boundingBox': '736,149,33,33', 'text': ' '}],

```

#### 4.1.1 Language and Orientation

Figure xx presents a sample that corresponds to the section 1 of the image containing Japanese text, but it also includes English text in some fragments of the text. The API offers the possibility to do self-detection, which in this case is helpful, but also the opportunity to specify the language in advance. Moreover, the orientation can also be self-detected or not.

#### 4.1.2 Regions

The blue line figure xx denotes the segment from the image that includes the characters. In this example, there is only one blue line because the document is only text. In other cases, it may contain images requiring more than one region.

一本の傘を二人でさす。多く男女二人の場合に  
いふ。あひがさ。

**アイアン** [Iron] (名) (普通にはなまつて「アイロ  
ン」といふ)。①鐵又は鐵製品。②西洋ひのし。③  
髪の毛をあらざる小形の鏡。④ゴルフ用の鐵の  
杖。—**アイシップリ** [Iron discipline] (名)  
【社】軍人に對する軍律、又は共產黨員に對する黨律  
の如く、絕對に破るべきからざる規律。—**ロー** [I-  
ron law] (名) ①人間の意思では變更することの  
できない法則。鐵則。②經賃銀鐵則。現在の資本主義  
經濟組織では、労働者の賃銀は増加させる望みがな  
いといふラッサールの説。即ち労働賃銀が上れば勞  
働者が増加して賃銀が低下する。故に労働者は生活  
費用として必要な額以上の賃銀を、永久得ることが  
出來ないとの説。

**アイ-レ-イ-ク** [愛育] (名) かはゆがつて育てる。こと。  
**アイ-レ-イ-シ-ヤ** [合醫者] (名) その人をよく診  
察し、通薦を授ける醫者。

**アイ-レ-イ-ン** [愛飲] (名) 好んで飲用すること。  
**アイ-レ-イ-ン** [合印] (名) 帳簿・書類のひきあはせ  
における印。わりいん。あひはん。

**アイヴァンホー** [Vanhoe] (名) 【文】イギリスの  
文豪サー・ウォルタースコットの歴史小説。騎士ア  
イヴァンホーを主人公としたもの。一八二〇年刊。

**アイヴオリ** [Ivory] (名) ①象牙。②名刺など  
に用ひる象牙色の光澤ある厚い洋紙。—**タワー**  
[Ivory tower] (名) ①象牙の塔。②現代物質文  
明の醜惡な生活を避け、自己の好む靜寂な詩的美の  
生活をするために文藝家などが一人で籠る別天地。

**アイ-エ-フ-テ-イ-ュ** [I-E-F-T-U] (名)  
【社】(International Federation of Trade Union)の  
略) 國際労働組合聯盟。第二インター・ナショナル系  
労働組合の國際的聯合で、赤色労働組合に對抗し、社  
會民主主義を標榜する。その本部がオランダのアム  
ステルダムにあるので、「アムステルダム・インター  
ナショナル」とも稱する。

**アイ-えん** [愛縁] (名) 【佛】恩愛の縁。

Figure 4: UPDATE CAPTION

一本の傘を二人でさす。多く男女二人の場合に  
いふ。あひがさ。

**アイアン** [Iron] (名) (普通にはなまつて「アイロ  
ン」といふ)。①鐵又は鐵製品。②西洋ひのし。③  
髪の毛をあらざる小形の鏡。④ゴルフ用の鐵の  
杖。—**アイシップリ** [Iron discipline] (名)  
【社】軍人に對する軍律、又は共產黨員に對する黨律  
の如く、絕對に破るべきからざる規律。—**ロー** [I-  
ron law] (名) ①人間の意思では變更することの  
できない法則。鐵則。②經賃銀鐵則。現在の資本主義  
經濟組織では、労働者の賃銀は増加させる望みがな  
いといふラッサールの説。即ち労働賃銀が上れば勞  
働者が増加して賃銀が低下する。故に労働者は生活  
費用として必要な額以上の賃銀を、永久得することが  
出來ないとの説。

**アイ-レ-イ-ク** [愛育] (名) かはゆがつて育てる。こと。  
**アイ-レ-イ-シ-ヤ** [合醫者] (名) その人をよく診  
察し、通薦を授ける醫者。

**アイ-レ-イ-ン** [愛飲] (名) 好んで飲用すること。  
**アイ-レ-イ-ン** [合印] (名) 帳簿・書類のひきあはせ  
における印。わりいん。あひはん。

**アイヴァンホー** [Vanhoe] (名) 【文】イギリスの  
文豪サー・ウォルタースコットの歴史小説。騎士ア  
イヴァンホーを主人公としたもの。一八二〇年刊。

**アイヴオリ** [Ivory] (名) ①象牙。②名刺など  
に用ひる象牙色の光澤ある厚い洋紙。—**タワー**  
[Ivory tower] (名) ①象牙の塔。②現代物質文  
明の醜惡な生活を避け、自己の好む靜寂な詩的美の  
生活をするために文藝家などが一人で籠る別天地。

**アイ-エ-フ-テ-イ-ュ** [I-E-F-T-U] (名)  
【社】(International Federation of Trade Union)の  
略) 國際労働組合聯盟。第二インター・ナショナル系  
労働組合の國際的聯合で、赤色労働組合に對抗し、社  
會民主主義を標榜する。その本部がオランダのアム  
ステルダムにあるので、「アムステルダム・インター  
ナショナル」とも稱する。

**アイ-えん** [愛縁] (名) 【佛】恩愛の縁。

Figure 5: UPDATE CAPTION

### 4.1.3 Lines

Once the region is defined, the API provides in figure xx (red rectangles) the subsets of characters of the region called Lines. From the example, we can see that some characters are not included in any of the lines. Those characters are not recognized by the API.

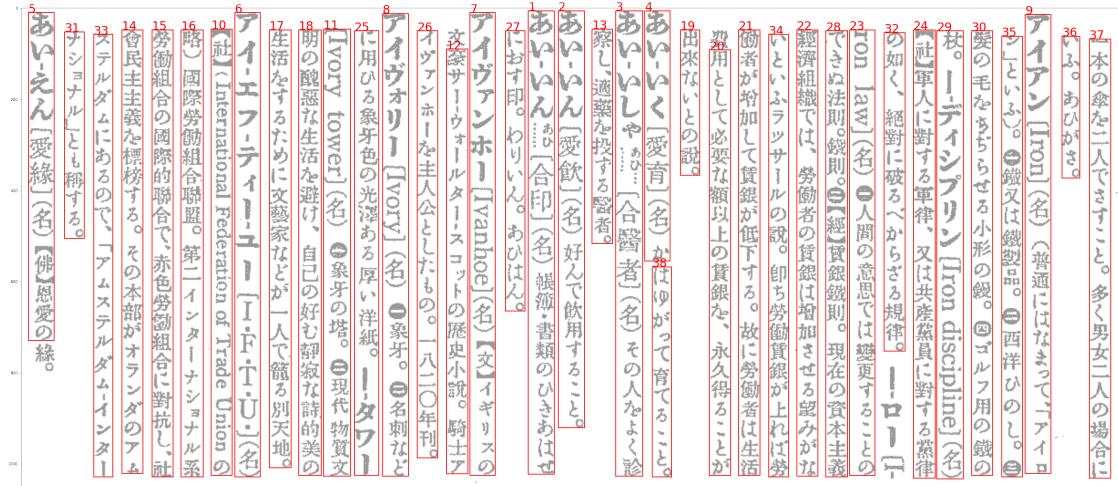


Figure 6: UPDATE CAPTION

### 4.1.4 The Character Boundaries

After the lines are identified, the API provides in figure xx each recognized character bounding boxes (yellow regions). Those without a yellow square were not recognized by the API in this example.



Figure 7: UPDATE CAPTION

### 4.1.5 Characters

Once the bounding boxes are located, each character takes its corresponding place, as it is shown in figure xx.

あいえん【愛縁】(名)【佛恩愛の縁。	一本の傘を三人でさすこと。多く男女二人の場合に いふ。あひがさ。
アイアン【iron】(名)【普通】(1)はなま。(2)アイロ ン」といふ。①鐵又は鐵製品。②西洋ひのし。③ 髮の毛をかららせる小形の鐵。④ゴルフ用の鐵の 枝。—・デイシブリン【Iron discipline】(名) 【社】軍人に對する軍律、又は共産黨員に對する黨律 の如く、絕對に破るべきからざる規律。—・ローラ ン law】(名)①人間の意思では變更することの できぬ法則。鐵則。②經濟鐵則。現在の資本主義 經濟組織では、労働者の賃銀は増加させる望みがな れといふラッサールの說。即ち労働賃銀が上れば労 働者が増加して賃銀が低下する。故に労働者は生活 費用として必要な額以上の賃銀を、永久得ることが 出来ないとの說。	一本の傘を二人でさすこと。多く男女二人の場合に いふ。あひがさ。
アイアンホー【Ivanhoe】(名)【文】イギリスの 文豪サー・ウォルタースコットの歴史小説。騎士ア イヴァンホーを主人公としたもの。一八二〇年刊。	一本の傘を二人でさすこと。多く男女二人の場合に いふ。あひがさ。
アイヴオリー【Ivory】(名)①象牙。②名刺など に用ひる象牙色の光澤ある厚い洋紙。—・タワー 【社】(International Federation of Trade Union Federation)国際労働組合聯盟。第一インター ナショナル系労働組合の國際労働組合で、赤色労働組合に對抗し、社會民主主義を標榜する。その本部がオランダのアム ステルダムにあるので、「アムステルダム・インター ナショナル」とも稱する。	一本の傘を二人でさすこと。多く男女二人の場合に いふ。あひがさ。

Figure 8: UPDATE CAPTION

#### 4.1.6 Advantages and Disadvantages

**4.1.6.1 Advantages** The API is easy to use after reading the documentation. It offers an example of how to use it in different programming languages. It also offers the possibility to self-detect languages and the orientation of the document. The response is clear, and extract the information from the JSON file is not complicated.

**4.1.6.2 Disadvantages** Even though the API offers the possibility to self-detect the language, It changes the order in which the characters should be in the output. In figure xx we can see that the first line provided by the self-detection approach is located in the middle. This situation does not allow a correct interpretation of the text because it does not follow the reading pattern. In this case, it is a vertical orientation, which means we have to read it from right to left.

あいえん【愛縁】(名)【佛恩愛の縁。	一本の傘を二人でさすこと。多く男女二人の場合に いふ。あひがさ。
アイアン【iron】(名)【普通】(1)はなま。(2)アイロ ン」といふ。①鐵又は鐵製品。②西洋ひのし。③ 髮の毛をかららせる小形の鐵。④ゴルフ用の鐵の 枝。—・デイシブリン【Iron discipline】(名) 【社】軍人に對する軍律、又は共産黨員に對する黨律 の如く、絕對に破るべきからざる規律。—・ローラ ン law】(名)①人間の意思では變更することの できぬ法則。鐵則。②經濟鐵則。現在の資本主義 經濟組織では、労働者の賃銀は増加させる望みがな れといふラッサールの說。即ち労働賃銀が上れば労 働者が増加して賃銀が低下する。故に労働者は生活 費用として必要な額以上の賃銀を、永久得することが 出来ないとの說。	一本の傘を二人でさすこと。多く男女二人の場合に いふ。あひがさ。
アイアンホー【Ivanhoe】(名)【文】イギリスの 文豪サー・ウォルタースコットの歴史小説。騎士ア イヴァンホーを主人公としたもの。一八二〇年刊。	一本の傘を二人でさすこと。多く男女二人の場合に いふ。あひがさ。
アイヴオリー【Ivory】(名)①象牙。②名刺など に用ひる象牙色の光澤ある厚い洋紙。—・タワー 【社】(International Federation of Trade Union Federation)国際労働組合聯盟。第一インター ナショナル系労働組合の國際労働組合で、赤色労働組合に對抗し、社會民主主義を標榜する。その本部がオランダのアム ステルダムにあるので、「アムステルダム・インター ナショナル」とも稱する。	一本の傘を二人でさすこと。多く男女二人の場合に いふ。あひがさ。

Figure 9: UPDATE CAPTION

When the parameter language is stated as Japanese since the beginning, as we can observe in figure xx, it recognizes the first line correctly.



Figure 10: UPDATE CAPTION

## 4.2 Google Vision OCR

Google Vision API can integrate vision detection features, including image labeling, face and landmark detection, optical character recognition (OCR), and tagging of explicit content. The text recognition process detects text in images and recognizes the text contained therein. Once detected, the recognizer then determines the actual text in each block and segments it into lines and words.

To evaluate this OCR option, we need the following python libraries.

```
from os import path
from PIL import Image
from PIL import ImageFont
from PIL import ImageDraw
from PIL import ImageEnhance
import base64 #Convert an image into base64 form for API requests
import requests #API request - GET/POST
import json
import ast
```

Same as the Azure API, we need to create the headers and parameters necessary for calling the API. We also need the python library `requests` to make the API call. The method `post` is needed as the API uses the HTTP method POST.

The following example illustrates how to state the headers and parameter, but also create two functions. The first one, called `encode_image` converts an image into base64 format, an encoding process designed to carry data stored in binary formats across channels that only reliably support text content. It can embed image files or other binary assets inside textual assets. The second one makes an HTTP POST request that calls Google's Vision API.

```
#Create the URL for the API call. The endpoint corresponds to the Azure's
# endpoint of the account.
```

```
url = "https://vision.googleapis.com/v1/images:annotate"

querystring = {"key": google_vision_api_key }
headers = {'Content-Type': "application/json"}

def encode_image(image_path):
    with open(image_path, "rb") as image_file:
        return base64.b64encode(image_file.read())

def image_request(image_path):

    payload = '{\"requests\": [{\"image\": {\"content\": \"' + \
    encode_image(image_path).decode('utf-8') + \
    '\"}, \"features\": [{\"type\": \"TEXT_DETECTION\"}]}]}'

    response = requests.request("POST", url, data=payload, headers=headers, params=querystring)

    return ast.literal_eval(response.text)
```

This is an example of the response:

```
{'responses': [{'textAnnotations': [{'locale': 'ja',  
    'description': '-]\n\n- [ ]()\n]() \n  0 \n    p®( ]() \n    \n \nV \n\n    \n',  
    'boundingPoly': {'vertices': [{'x': 26, 'y': 21},  
        {'x': 1556, 'y': 21},  
        {'x': 1556, 'y': 2748},  
        {'x': 26, 'y': 2748}]},  
    {'description': '-',  
    'boundingPoly': {'vertices': [{'x': 1227, 'y': 686},  
        {'x': 1227, 'y': 677},  
        {'x': 1262, 'y': 677},  
        {'x': 1262, 'y': 686}]},  
    '\n'}]}]  
    {'description': '',  
    'boundingPoly': {'vertices': [{'x': 785, 'y': 880},  
        {'x': 814, 'y': 880},  
        {'x': 814, 'y': 914},  
        {'x': 785, 'y': 914}]},  
    {'description': '(',  
    'boundingPoly': {'vertices': [{'x': 775, 'y': 916},  
        {'x': 814, 'y': 916},  
        {'x': 814, 'y': 940},  
        {'x': 775, 'y': 940}]},  
    ...],  
    ....... ) \n-[ ~Q)
```

From the previous example, we can observe that the response returns a list of nested dictionaries. The first nested dictionary called **textAnnotations** consists of a list of dictionaries. Each one of those has a description (character) and boundingPoly (boundaries of the characters).

### 4.2.1 BoundingPoly

Each boundingPoly consist of vertices in which the symbol is located.

**4.2.1.1 Region** The first object in the nested dictionary contains all identified characters in the image as well as the boundaries that comprises the location of those characters. This object can resemble the idea of the region provided in the Azure's API.

Figure xx shows the region identified by the API.

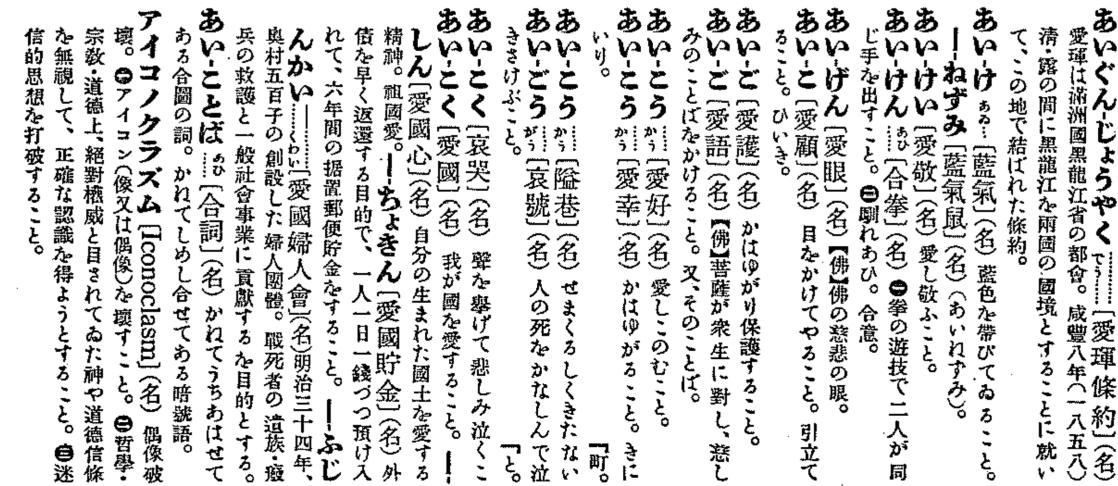


Figure 11: UPDATE CAPTION

**4.2.1.2 The Character Boundaries** The remaining objects enclose the boundaries of the characters. As it is shown in figure xx, each boundingPoly node can include one or more characters.

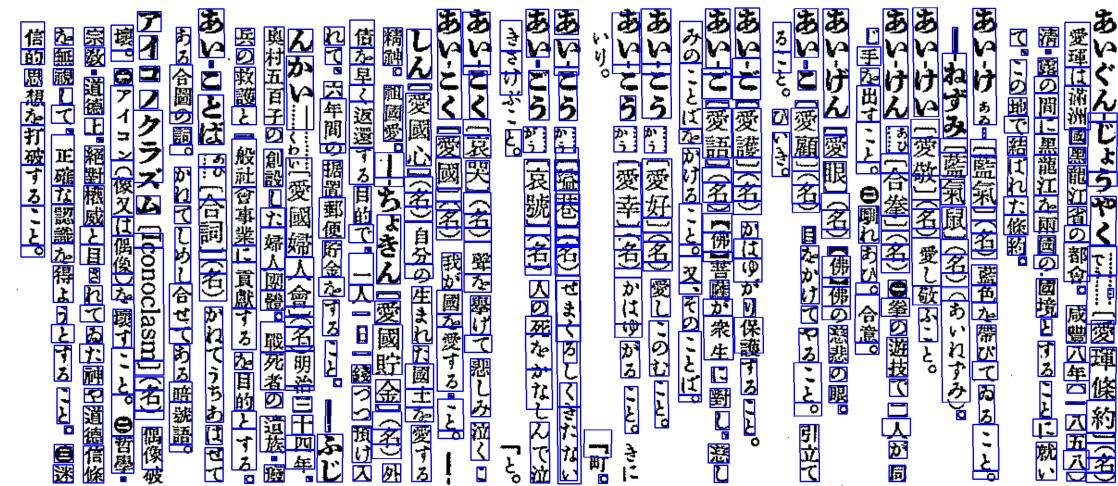


Figure 12: UPDATE CAPTION

### 4.2.2 Characters

Once the bounding boxes are located, each character takes its corresponding place, as it is shown in figure xx. Pending find a way to plot the text in vertical form to make it readable.



Figure 13: UPDATE CAPTION

### 4.2.3 languageHints

Google Vision API offers the advantage of including in the request language hints. It consists of a list of languages the user knows in advance the image contains. For this study, we run two versions, one without hints and the other one with a list of two hints, such as Japanese and English, as the image contains some words in English. Both versions returned the same output.

#### 4.2.4 Advantages and Disadvantages

**4.2.4.1 Advantages** Same as the Azure API, this is also easy to use after reading the documentation. It also offers several examples in different programming languages using the Google libraries for Python. It gives the possibility to include languages in advanced as hints, and identifies the orientation automatically. The response contains an object with all identified characters in the correct order.

Moreover, the API gives the opportunity to send the image in base64 format or using an internal (google storage) or external URL.

**4.2.4.2 Disadvantages** Unlike Azure Cognitive Services, which returns each recognized character individually, Google Vision groups characters in some cases. This situation adds some degree of complexity when processing the response from the API. [Not sure if it is a disadvantage in real life]

### 4.3 PyTesseract - Python Library

Python-Tesseract is an optical character recognition (OCR) library for python. It is a wrapper for Google's Tesseract-OCR Engine. PyTesseract can read all image types such as png, jpeg, gif, tiff, BMP, among others. It is widely used to process everything from scanned documents. This inbuilt python library performs the OCR techniques on the image to produce the digitized output text in a readable and editable form.

To use this library, we need to install Tesseract. For this study, we use version 4. It is necessary to call the .exe file to use Tesseract with PyTesseract. See the following example.

```
#Call .exe location for Windows  
pytesseract.pytesseract.tesseract_cmd = r'C:\Program Files\Tesseract-OCR\tesseract.exe'
```

Along with the PyTesseract library, the following libraries are used.

```
from PIL import Image  
from PIL import ImageOps
```

```

import pytesseract
from os import path
import pandas as pd
import matplotlib.pyplot as plt
from matplotlib.patches import Rectangle

```

PyTesseract offers several methods with different outputs when processing an image to extract text. According to Pypi.org, the uses of those methods are:

- **image\_to\_string**: Returns the result of a Tesseract OCR run on the image to string
- **image\_to\_boxes** Returns result containing recognized characters and their box boundaries
- **image\_to\_data**: Returns result containing box boundaries, confidences, and other information. Requires Tesseract 3.05+. For more information, please check the Tesseract TSV documentation
- **image\_to\_osd**: Returns result containing information about orientation and script detection.
- **image\_to\_alto\_xml**: Returns result in the form of Tesseract's ALTO XML format.
- **run\_and\_get\_output**: Returns the raw output from Tesseract OCR. Gives a bit more control over the parameters that are sent to tesseract.

For this analysis, we use the method `image_to_data` as this offers the boundaries and the recognized characters in a handy form so we can understand the process. Furthermore, in PyTesseract, it is needed to specify the language in advance. Since PyTesseract is a wrapper of Google's Tesseract-OCR Engine, Tesseract-OCR must have installed all required languages. In this case, it is Japanese, which is not included in its out-of-the-box version. To include a new language, adding the `.traineddata` file into the Tesseract-OCR/tessdata folder is sufficient. However, since the text in the image is vertical Japanese. It is necessary to include two files: `jpn.traineddata` and `jpn_vert.traineddata`. Once the files are located correctly, the following code extracts the text from the image.

```

img1 = Image.open(r'~\inputs\00-sample-pages_1.png')
image_to_data = pytesseract.image_to_data(img1, lang = 'jpn_vert')

```

The output from the previous code converted into a pandas dataframe is shown in figure xx.

	1	result.head(8)
	1	level page_num block_num par_num line_num word_num left top width height conf text
0	1	1 0 0 0 0 0 0 0 0 1588 2776 -1 NaN
1	2	1 1 0 0 0 0 898 30 654 657 -1 NaN
2	3	1 1 1 1 0 0 898 30 654 657 -1 NaN
3	4	1 1 1 1 1 0 1522 70 30 616 -1 NaN
4	5	1 1 1 1 1 1 1522 70 28 44 13 本
5	5	1 1 1 1 1 2 1525 123 26 35 92 の
6	5	1 1 1 1 1 3 1525 170 25 20 87 翠
7	5	1 1 1 1 1 4 1525 199 26 20 78 を

Figure 14: UPDATE CAPTION

From figure xx, we can observe that the output in the column **level** generates five levels. Those are pages, blocks, paragraphs, lines, and words (characters). Each one is a subset of the previous one.

### 4.3.1 Page

For this example, the page resembles the concept of region, which is the largest area of the image enclosing all recognized characters. Figure 12 illustrates the page area.

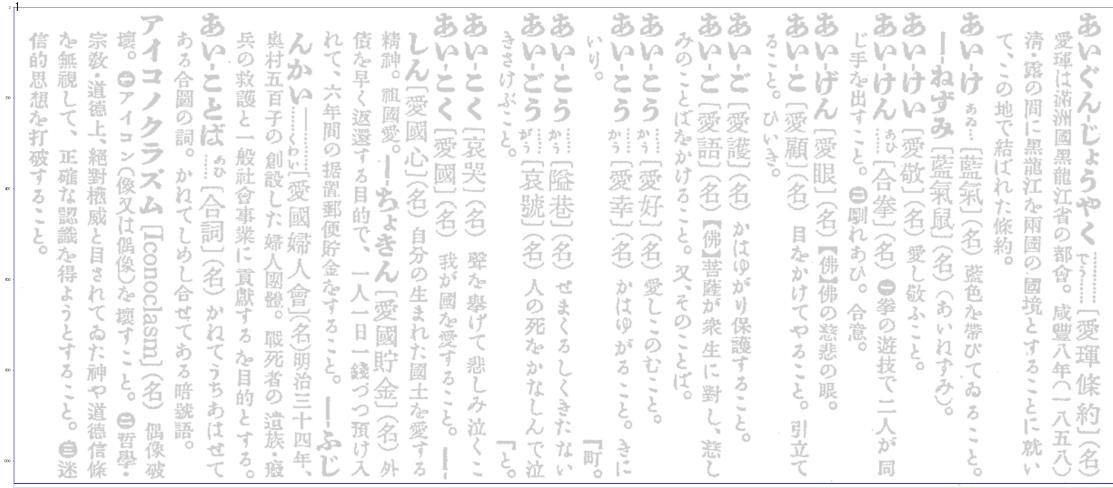


Figure 15: UPDATE CAPTION

### 4.3.2 Blocks

The blocks are subsets of the page area that consist of small groups of recognized characters. In figure 13, we can see that most of the blocks are well ordered. In terms of correctness of the boundaries, we can see that some of them overlap several characters that are not recognized by the library.

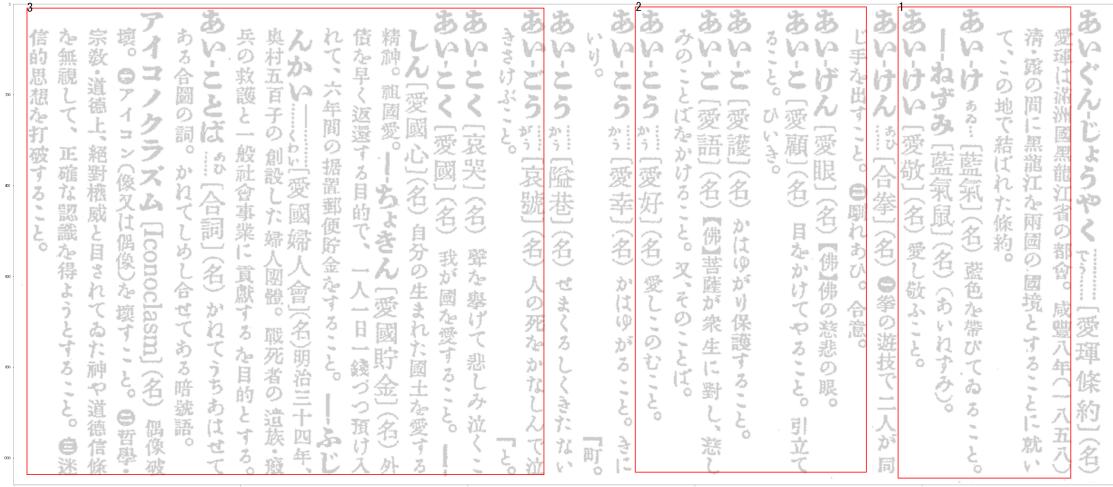


Figure 16: UPDATE CAPTION

### 4.3.3 Paragraphs

Similar to the blocks, the paragraphs are subsets of blocks. These follow the same idea of subareas containing a set of characters. In this case, for example, we can look at paragraphs three and four, which are subsets of block four (figure xx). It is interesting to see how the only paragraph four of block four is a very small area compared to its parent block. This situation happens in several other blocks. See figure xx for a better understanding.

32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	
信的思惟を打破すること。	32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1

あいぐんじょうやく [愛珲條約] (名)  
愛珲は瀋洲國黒龍江省の都會。咸豐八年(一八五八)清・露の間に黒龍江を兩國の國境とすることに就いて、この地で結ばれた條約。

あいねずみ [藍氣鼠] (名) 藍色を帶びてゐること。  
あいねずみ [藍氣鼠] (名) (あいねずみ)。

あいけい [愛敬] (名) 愛し敬ふこと。  
あいけい [愛敬] (名) (あいけい)。

あいけん [合拳] (名) ①拳の遊技で二人が同じ手を出すこと。②剛れおひ。合意。

あいこ [愛顧] (名) 目をかけてやること。引立てること。  
あいこ [愛顧] (名) (ひいき)。

あいご [愛護] (名) かはゆがり保護すること。  
あいご [愛護] (名) (かはゆがり)。

あいご [愛護] (名) 善薩が衆生に對し、慈しむこと。  
あいご [愛護] (名) (かはゆがり)。

あいこう [愛好] (名) 愛しきのむこと。  
あいこう [愛好] (名) (かはゆがる)。

あいこう [愛幸] (名) かはゆがるのこと。  
あいこう [愛幸] (名) (かはゆがる)。

あいこう [隘巷] (名) せまくるしくきたない町。  
あいこう [隘巷] (名) (かいきょう)。

あいごう [哀號] (名) 人の死をかなしんで泣き声をあげること。  
あいごう [哀號] (名) (あいごう)。

あいごく [哀哭] (名) 翁を擧げて悲しみ泣くこと。  
あいごく [哀哭] (名) (あいごく)。

あいごく [愛國] (名) 我が國を愛すること。  
あいごく [愛國] (名) (あいごく)。

あいごく [愛國心] (名) 自分の生まれた國土を愛する精神。祖國愛。  
あいごく [愛國心] (名) (あいごく)。

あいこう [愛國貯金] (名) 外債を早く返還する目的で、一人一日一錢づつ預け入れて、六年間の据置郵便貯金をすること。  
あいこう [愛國貯金] (名) (あいこう)。

あいごと [愛國] (名) 我が國を愛すること。  
あいごと [愛國] (名) (あいごと)。

あいごと [愛護] (名) かはゆがり保護すること。  
あいごと [愛護] (名) (かはゆがり)。

あいごと [愛語] (名) 【佛】菩薩が衆生に對し、慈しむこと。  
あいごと [愛語] (名) (かはゆがり)。

あいごと [愛好] (名) 愛しきのむこと。  
あいごと [愛好] (名) (かはゆがる)。

あいこう [愛幸] (名) かはゆがるのこと。  
あいこう [愛幸] (名) (かはゆがる)。

あいこう [隘巷] (名) せまくるしくきたない町。  
あいこう [隘巷] (名) (かいきょう)。

あいこう [愛國] (名) 我が國を愛すること。  
あいこう [愛國] (名) (あいこう)。

あいこう [愛國心] (名) 自分の生まれた國土を愛する精神。祖國愛。  
あいこう [愛國心] (名) (あいこう)。

あいこう [愛國貯金] (名) 外債を早く返還する目的で、一人一日一錢づつ預け入れて、六年間の据置郵便貯金をすること。  
あいこう [愛國貯金] (名) (あいこう)。

あいごと [愛國] (名) 我が國を愛すること。  
あいごと [愛國] (名) (あいごと)。

あいごと [愛護] (名) かはゆがり保護すること。  
あいごと [愛護] (名) (かはゆがり)。

あいごと [愛語] (名) 【佛】菩薩が衆生に對し、慈しむこと。  
あいごと [愛語] (名) (かはゆがり)。

あいごと [愛好] (名) 愛しきのむこと。  
あいごと [愛好] (名) (かはゆがる)。

あいこう [愛幸] (名) かはゆがるのこと。  
あいこう [愛幸] (名) (かはゆがる)。

Figure 17: UPDATE CAPTION

#### 4.3.4 Lines

Once the paragraphs are defined, the library provides in figure xx, the last subsets of characters called Lines. From the example, we can see that most of the characters are located in one of the lines. However, several others are not. For some cases there is no clear reason why those characters are not included. For instance, between the lines 14 and 15, two clear lines are not recognized by the library.

32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	
信的思惟を打破すること。	32	31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1

あいぐんじょうやく [愛珲條約] (名)  
愛珲は瀋洲國黒龍江省の都會。咸豐八年(一八五八)清・露の間に黒龍江を兩國の國境とすることに就いて、この地で結ばれた條約。

あいねずみ [藍氣鼠] (名) 藍色を帶びてゐること。  
あいねずみ [藍氣鼠] (名) (あいねずみ)。

あいけい [愛敬] (名) 愛し敬ふこと。  
あいけい [愛敬] (名) (あいけい)。

あいけん [合拳] (名) ①拳の遊技で二人が同じ手を出すこと。②剛れおひ。合意。

あいこ [愛顧] (名) 目をかけてやること。引立てること。  
あいこ [愛顧] (名) (ひいき)。

あいご [愛護] (名) かはゆがり保護すること。  
あいご [愛護] (名) (かはゆがり)。

あいご [愛護] (名) 善薩が衆生に對し、慈しむこと。  
あいご [愛護] (名) (かはゆがり)。

あいこう [愛好] (名) 愛しきのむこと。  
あいこう [愛好] (名) (かはゆがる)。

あいこう [愛幸] (名) かはゆがるのこと。  
あいこう [愛幸] (名) (かはゆがる)。

あいこう [隘巷] (名) せまくるしくきたない町。  
あいこう [隘巷] (名) (かいきょう)。

あいこう [哀號] (名) 人の死をかなしんで泣き声をあげること。  
あいこう [哀號] (名) (あいこう)。

あいごく [哀哭] (名) 翁を擧げて悲しみ泣くこと。  
あいごく [哀哭] (名) (あいごく)。

あいごく [愛國] (名) 我が國を愛すること。  
あいごく [愛國] (名) (あいごく)。

あいごと [愛護] (名) かはゆがり保護すること。  
あいごと [愛護] (名) (かはゆがり)。

あいごと [愛語] (名) 【佛】菩薩が衆生に對し、慈しむこと。  
あいごと [愛語] (名) (かはゆがり)。

あいごと [愛好] (名) 愛しきのむこと。  
あいごと [愛好] (名) (かはゆがる)。

あいこう [愛幸] (名) かはゆがるのこと。  
あいこう [愛幸] (名) (かはゆがる)。

Figure 18: UPDATE CAPTION

#### 4.3.5 Character Boundaries

These boundaries enclose the characters. Figure xx shows that a character boundary can include one or more characters. However, some character boundaries overlap each other. Moreover, same as the previous areas, several characters have been skipped with no apparent reason.

#### 4.3.6 Characters

In terms of character-level recognition, we can observe that many characters are not recognized. Interestingly those characters are in a correct block, paragraph, line, and character boundary, but it was skipped. For

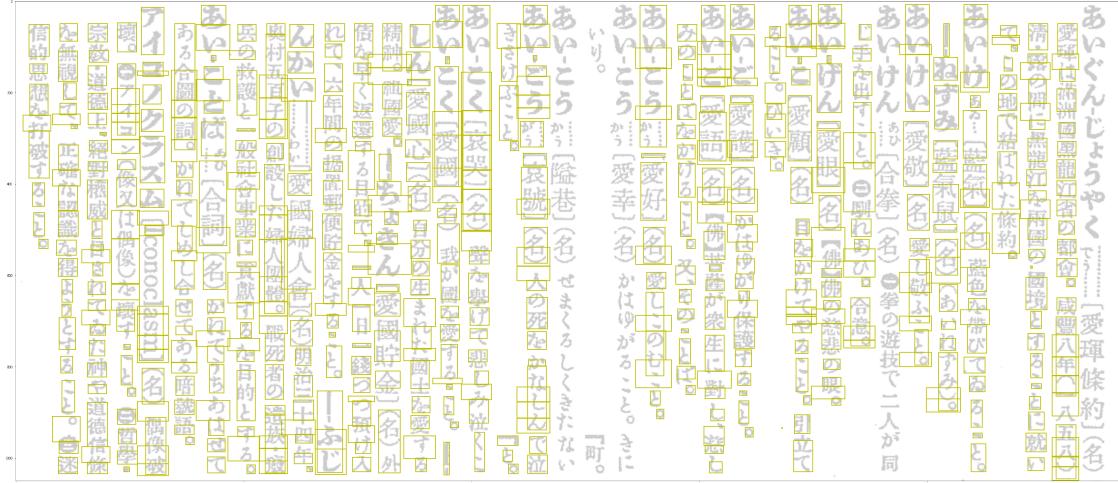


Figure 19: UPDATE CAPTION

example, in the first line of recognized characters and the first character boundary of this line, there are two words, but one of them is not recognized. Figure xx illustrates this case.

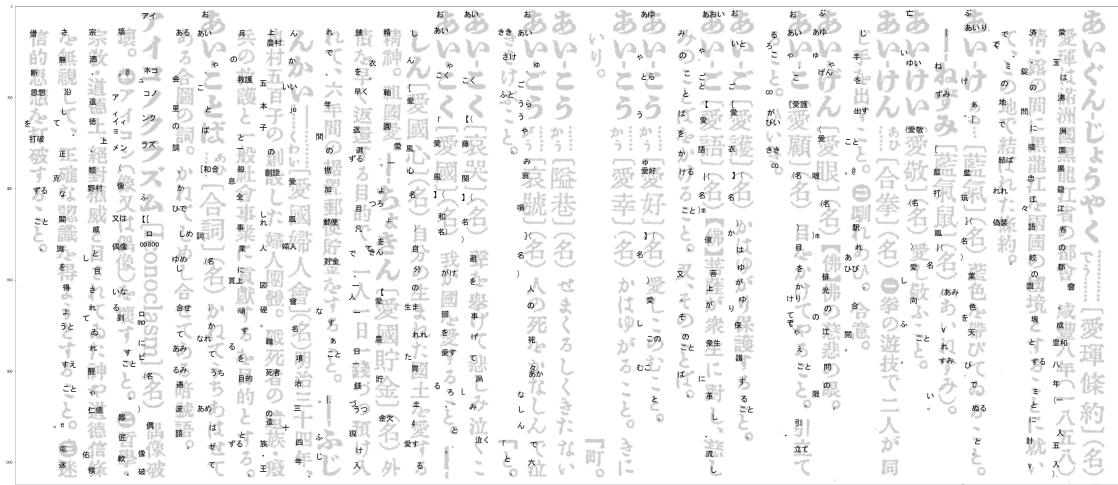


Figure 20: UPDATE CAPTION

### 4.3.7 Advantages and Disadvantages

**4.3.7.1 Advantages** The most important advantage is that this is a free library. It is also a wrapper from Google's Tesseract Engine. Moreover, the way it presents the covered areas are very detailed. It offers blocks, paragraphs, lines, boundaries, and finally the character level.

**4.3.7.2 Disadvantages** Some lines were completely ignored by Pyteseract, even thought it was clear enough to be recognized by the tool. Moreover, same as Google Vision, this library groups more than one character sometimes.

## 4.4 Kindai-OCR

According to the Center for Open Data in the Humanities, Kindai-OCR is an OCR system for modern Japanese documents. It was built under the N2I project that study state-of-the-art statistical models for

OCR, and work on the development of infrastructure and providing open access to data. The datasets were constructed by the University of Tokyo Center for Education and Research and the National Institute of Japanese Language and are used for OCR machine learning. These datasets are commonly known as the Modern Magazine Datasets.

The Kindai-OCR code was developed in Python and released as an open source code in August 2020 on Github. The Github folder contains all the necessary code and instructions to assess the OCR tool. [not sure if it is necessary to mention the modification in the file test.py I did, in order to run the model using GPUs.]

After running the Kindai-OCR for our samples, the system generated two main outputs. They are explained as follows:

#### 4.4.1 File result.xml

This file is presented in an Extensible Markup Language (XML) format containing a tree with the identified lines, as well as the corresponding characters. An example of the xml file is shown as follows:

```
<?xml version='1.0' encoding='Shift_JIS'?>
<paper xmlns="http://codh.rois.ac.jp/modern-magazine/"><page dpi="100" file="00-sample_pages_3_sec_3.png">
```

After transforming this file into a dataframe form, we found that 238 lines with their corresponding characters, were generated. In figure XX we can see that the result.xml file provides information such as the corresponding level (paper, page, or line) in the tag column, the identified characters in the text column, the associated image the file column, among other attributes.

	tag	attributes	text	dpi	file	number	width	height	x	y
0	{http://codh.rois.ac.jp/modern-magazine/paper}	{} None								
1	{http://codh.rois.ac.jp/modern-magazine/}page	{"dpi": "100", "file": "00-sample_pages_3_sec_3.png"} None		100	sample_pages_3_sec_3.png	1	2430	1040		
2	{http://codh.rois.ac.jp/modern-magazine/}line	{"height": "905", "width": "72", "x": "あいくるし發路事。愛くるし (形。..."} 100	あいくるし發路事。愛くるし (形。	100	sample_pages_3_sec_3.png	1	72	905	332	0
3	{http://codh.rois.ac.jp/modern-magazine/}line	{"height": "1026", "width": "76", "x": "あいくち島 (合巴, 名。あれせも。ものよ..."} 100	あいくち島 (合巴, 名。あれせも。ものよ	100	sample_pages_3_sec_3.png	1	76	1026	466	0
4	{http://codh.rois.ac.jp/modern-magazine/}line	{"height": "1026", "width": "69", "x": "あいくすりし ((きき)名..."} 100	あいくすりし ((きき)名その人の萬に信身あ	100	sample_pages_3_sec_3.png	1	69	1026	592	0

Figure 21: UPDATE CAPTION

#### 4.4.2 New version of the image

Kindai-OCR generated a new version of the image containing the identified lines and characters. An example is shown as follows in figure xx.

#### 4.4.3 Advantages and Disadvantages

**4.4.3.1 Advantages** Same as libraries such as Pytesseract, one of The most important advantages is that it is a free tool. Moreover, the recognition process is very easy. To obtain results, just by putting the images in the test folder and running the file called test.py on Python is enough to obtain the results.

**4.4.3.2 Disadvantages** At the time of this assessment, the use of GPUs is required which is a limitation to take into account. For this study, we used Colab from Google, which offers this capability. Moreover, Kindai-OCR returns the characters grouped by lines. They are not generated individually.

## 4.5 Kraken - Python Library

Kraken is a python library created by Benjamin Kiessling from Université PSL in France and Leipzig University. This library is a combination of a fork (independent development built on an existing one) from the

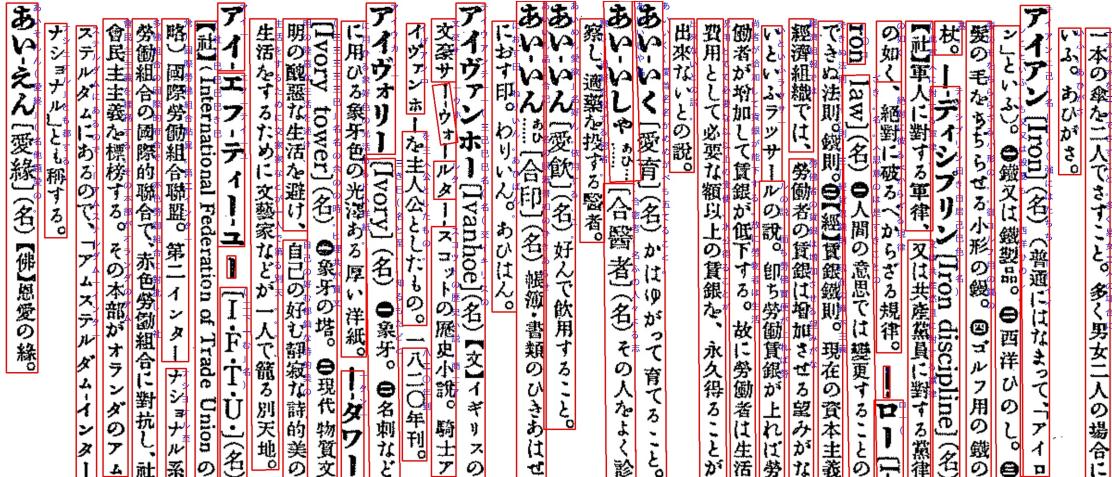


Figure 22: UPDATE CAPTION

Ocropus package and CLSTM neural network library. Its purpose was to improve accuracy and performance rates compared with Ocropolis. It has been tested with languages such as Arabic, Persian, Syriac Polytonic Greek, Latin, among other. However, there is no information about japanese-trained models.

To assess the performance of this library, it is necessary to train a model for japanese texts. To do that, Kraken offers a function called **train**. This function expects a series of images that according to Kraken's website, these images must be high quality scans, preferably color or grayscale, and at least 300dpi. Images in PDF format are not allowed. Each image must be accompanied with a text file with extension .gt.txt. This file will contain the character(s) that correspond(s) to the image with the same name. For instance in figure xx we see the image file is called **sec1\_line1.png**, the pair-wise text element for this image is **sec1\_line1.gt.txt**.

After the training stage, Kraken analyses the layout and extracts text lines from an input image for later processing by the recognition. A step by step process is provided by Kraken's website (<http://kraken.re/index.html>) to recognize text from an image using the model we want.

#### 4.5.1 Training process challenges

The training process itself is very friendly. However, the data preparation is the most important part as it generates the samples we feed into the model. Other studies pointed out that for languages such as Arabic, at least 800 lines are required to train a descent model. For the Japanese text of this study we trained 36 lines in which several characters appear more than once. The result after the training process was not good, indicating that the need of greater sample data is needed, which is not in the scope of this study.

#### 4.5.2 Other limitations

In order to use this library it is required to have either Linux or Mac. In this study, we use Google Colab which offers a Linux environment to assess the library.

### 4.6 Amazon Textract

According to Amazon's website, Amazon Textract can detect Latin-script characters from the standard English alphabet and ASCII symbols.

<https://aws.amazon.com/textract/faqs/>

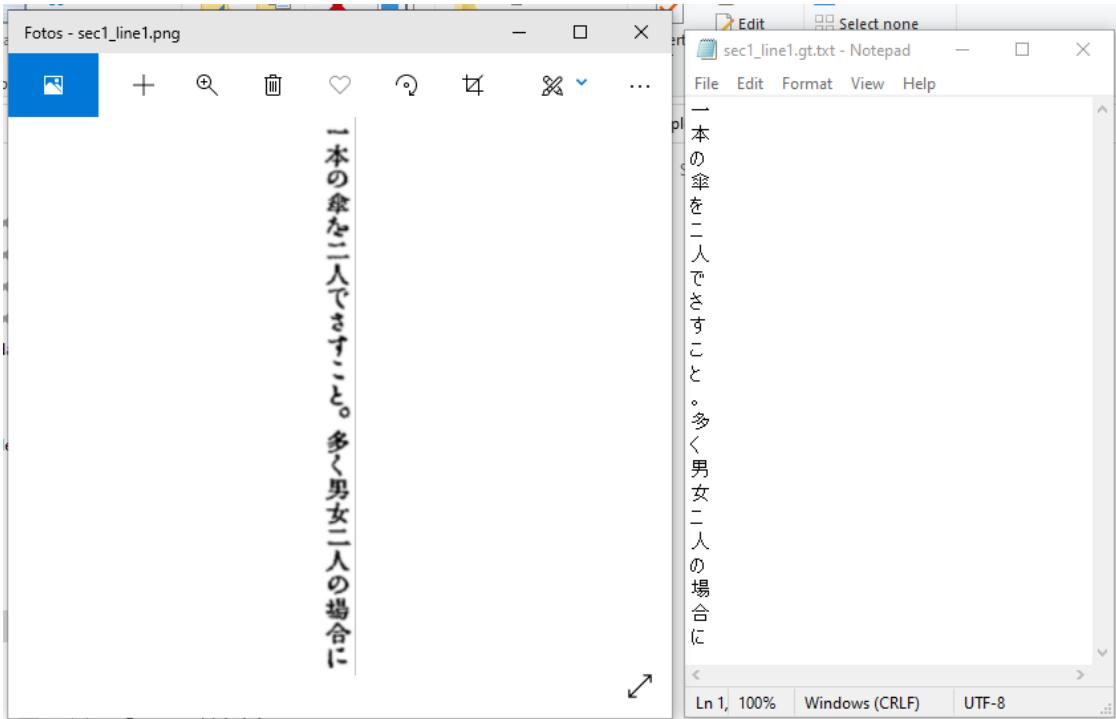


Figure 23: UPDATE CAPTION

## 5 Evaluation

For each option, we rank them based on two scales: results and difficulty. We need to work out a definition for each of these.

Diego's comment: From what I learned so far, it could be easiness (an API is more straightforward than a library), number of characters correctly recognized, flexibility, and cost. Moreover, some APIs such as Google's groups more than one character, whereas Azure returns only one character per boundary.

In terms of flexibility...

Sequence of the characters...

Output given...xml, json, flat file, csv...

Number of characters well recognized...(Accuracy)

The main aspect of the end product will be a graph with two axis, and dots for where each service is positioned, coloured or faceted by the test.

## 6 Discussion

## 7 Appendix

1: [https://www.researchgate.net/publication/267465115\\_An\\_Overview\\_and\\_Applications\\_of\\_Optical\\_Character\\_Recognition](https://www.researchgate.net/publication/267465115_An_Overview_and_Applications_of_Optical_Character_Recognition)

2: <https://books-scholarsportal-info.myaccess.library.utoronto.ca/en/read?id=/ebooks/ebooks2/wiley/2011-12-13/1/9780470176535>

4: <https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/home>

- 5: <https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/quickstarts/python-print-text>
- 6: <https://cloud.google.com/vision/docs/ocr>
- 7: <https://pypi.org/project/pytesseract/>
- 8: <https://github.com/tesseract-ocr/tessdata>
- 9: <https://github.com/UB-Mannheim/tesseract/wiki>

## References

- “Amazon Textract FAQs.” 2020. <https://aws.amazon.com/textract/faqs/>.
- Bowley, Thomas. 1902. *Elements of Statistics*.
- Breuel, Thomas M. 2007. “Announcing the OCropus Open Source OCR System.” <https://developers.googleblog.com/2007/04/announcing-ocropus-open-source-ocr.html>.
- Brontë, Charlotte. 1847. *Jane Eyre*.
- Cheriet, Mohamed, Nawwaf Kharma, Cheng-Lin Liu, and Ching Y. Suen. 2007. *Character Recognition Systems: A Guide for Students and Practitioner*. Wiley.
- “Google Vision API Documentation.” 2020. <https://cloud.google.com/vision>.
- Horiuchi, Tadashi, and Satoru Kato. 2009. “A Study on Japanese Historical Character Recognition Using Modular Neural Networks.” In *2009 Fourth International Conference on Innovative Computing, Information and Control (Icicic)*, 1507–10. <https://doi.org/10.1109/ICICIC.2009.57>.
- Kanagarathinam, Karthick, K Ravindrakumar, R Francis, and S Ilankannan. 2019. “Steps Involved in Text Recognition and Recent Research in Ocr; a Study” 8 (May): 3095–3100.
- Kiessling, Benjamin. 2019. “Kraken - an Universal Text Recognizer for the Humanities.” <https://dev.clariah.nl/files/dh2019/boa/0673.html>.
- Kim, Min-Soo, Man-Dae Jang, Hyun-Il Choi, Taik-Heon Rhee, Jin-Hyung Kim, and Hee-Kue Kwag. 2004. “Digitalizing Scheme of Handwritten Hanja Historical Documents.” In *First International Workshop on Document Image Analysis for Libraries, 2004. Proceedings*, 321–27. IEEE.
- Le, Anh Duc, Tarin Clanuwat, and Asanobu Kitamoto. 2019. “A Human-Inspired Recognition System for Pre-Modern Japanese Historical Documents.” *IEEE Access* 7: 84163–9. <https://doi.org/10.1109/ACCESS.2019.2924449>.
- Le, Anh Duc, Daichi Mochihashi, Katsuya Masuda, Hideki Mima, and Nam Tuan Ly. 2019. “Recognition of Japanese Historical Text Lines by an Attention-Based Encoder-Decoder and Text Line Generation.” In *Proceedings of the 5th International Workshop on Historical Document Imaging and Processing*, 37–41.
- Lombardi, Francesco, and Simone Marinai. 2020. “Deep Learning for Historical Document Analysis and Recognition—a Survey.” *Journal of Imaging* 6 (10): 110.
- Losee, Robert M. 2006. “Browsing Mixed Structured and Unstructured Data.” *Information Processing & Management* 42 (2): 440–52.
- Martínek, Jiří, Ladislav Lenc, and Pavel Král. 2020. “Building an Efficient Ocr System for Historical Documents with Little Training Data.” *Neural Computing and Applications* 31: 17209–27. <https://doi.org/10.1007/s00521-020-04910-x>.
- “Microsoft Azure Documentation.” 2020. <https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/concept-recognizing-text#supported-languages-for-print-text>.
- Nguyen, Hung Tuan, Nam Tuan Ly, Kha Cong Nguyen, Cuong Tuan Nguyen, and Masaki Nakagawa. 2017. “Attempts to Recognize Anomalously Deformed Kana in Japanese Historical Documents.” In *Proceedings of the 4th International Workshop on Historical Document Imaging and Processing*, 31–36. HIP2017. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3151509.3151514>.
- Ooms, Jeroen. 2019. *Tesseract: Open Source Ocr Engine*. <https://CRAN.R-project.org/package=tesseract>.
- Reul, Christian, Dennis Christ, Alexander Hartelt, Nico Balbach, Maximilian Wehner, Uwe Springmann, Christoph Wick, Christine Grundig, Andreas Büttner, and Frank Puppe. 2019. “OCR4all—an Open-Source Tool Providing a (Semi-) Automatic Ocr Workflow for Historical Printings.” *Applied Sciences* 9 (22): 4853.
- Rodrigues, Bruno. 2019. “Historical newspaper scraping with tesseract and R.” [https://www.brodrigues.co/blog/2019-04-07-historical\\_newspaper\\_scraping\\_tesseract/#](https://www.brodrigues.co/blog/2019-04-07-historical_newspaper_scraping_tesseract/#).

Romanov, Maxim, Matthew Thomas Miller, Sarah Bowen Savant, and Benjamin Kiessling. 2017. “Important New Developments in Arabographic Optical Character Recognition (Ocr).” *arXiv Preprint arXiv:1703.09550*.

Rychlik, Marek, Dwight Nwaigwe, Yan Han, and Dylan Murphy. 2020. “Development of a New Image-to-Text Conversion System for Pashto, Farsi and Traditional Chinese.” <http://arxiv.org/abs/2005.08650>.

Smith, Ray W. 2013. “History of the Tesseract OCR engine: what worked and what didn’t.” In *Document Recognition and Retrieval Xx*, edited by Richard Zanibbi and Bertrand Coësnon, 8658:1–12. International Society for Optics; Photonics; SPIE. <https://doi.org/10.1117/12.2010051>.

“Tesseract Manual Page.” 2018. <https://github.com/tesseract-ocr/tesseract/blob/master/doc/tesseract.1.asc>.

Valy, D., M. Verleysen, and S. Chhun. 2020. “Data Augmentation and Text Recognition on Khmer Historical Manuscripts.” In *2020 17th International Conference on Frontiers in Handwriting Recognition (Icfhr)*, 73–78. <https://doi.org/10.1109/ICFHR2020.2020.00024>.