

OPTICAL CHARACTER RECOGNITION (OCR) WITH ABBYY FINEREADER & KRAKEN

Institute for Computational Linguistics, 2020

Abstract: This paper briefly summarizes the applied OCR methods as well as the (optional) file splitting in the digitalization process of 10'012 file cards, given as pages in one large pdf file (778.8 MB).

1 Preprocessing

Breaking the input into multiple pieces is not necessarily needed, however, «divide et impera» appears to be a recommendable strategy: A preceding splitting beware from reprocessing everything in case of an error, and might, facilitating parallelization, speed up the process significantly.

```
1 from PyPDF2 import PdfFileWriter, PdfFileReader
2 input_pdf = PdfFileReader(open("/path/to/input/file.pdf", "rb"))
3 for i in range(input_pdf.numPages):
4     with open("file_"+(5-len(str(i+1)))*'0'+str(i+1)+".pdf", "wb") as output:
5         PdfFileWriter().addPage(input_pdf.getPage(i)).write(output)
```

Code Snipped: A Python-script to save each page of a pdf file as a new pdf file.

The names of the output files should be of uniform length (depending on the particular naming system). Otherwise, *abc*-sort of certain file explorers might destroy the order, which eventually has disadvantageous effects on further processing steps (e.g. required id-extraction from file names).

2 ABBYY FineReader

This section documents the process of applying OCR on multiple files by the commercially proprietary software «FineReader» by ABBYY, which involves the creation of different directories as well as some configuration steps.

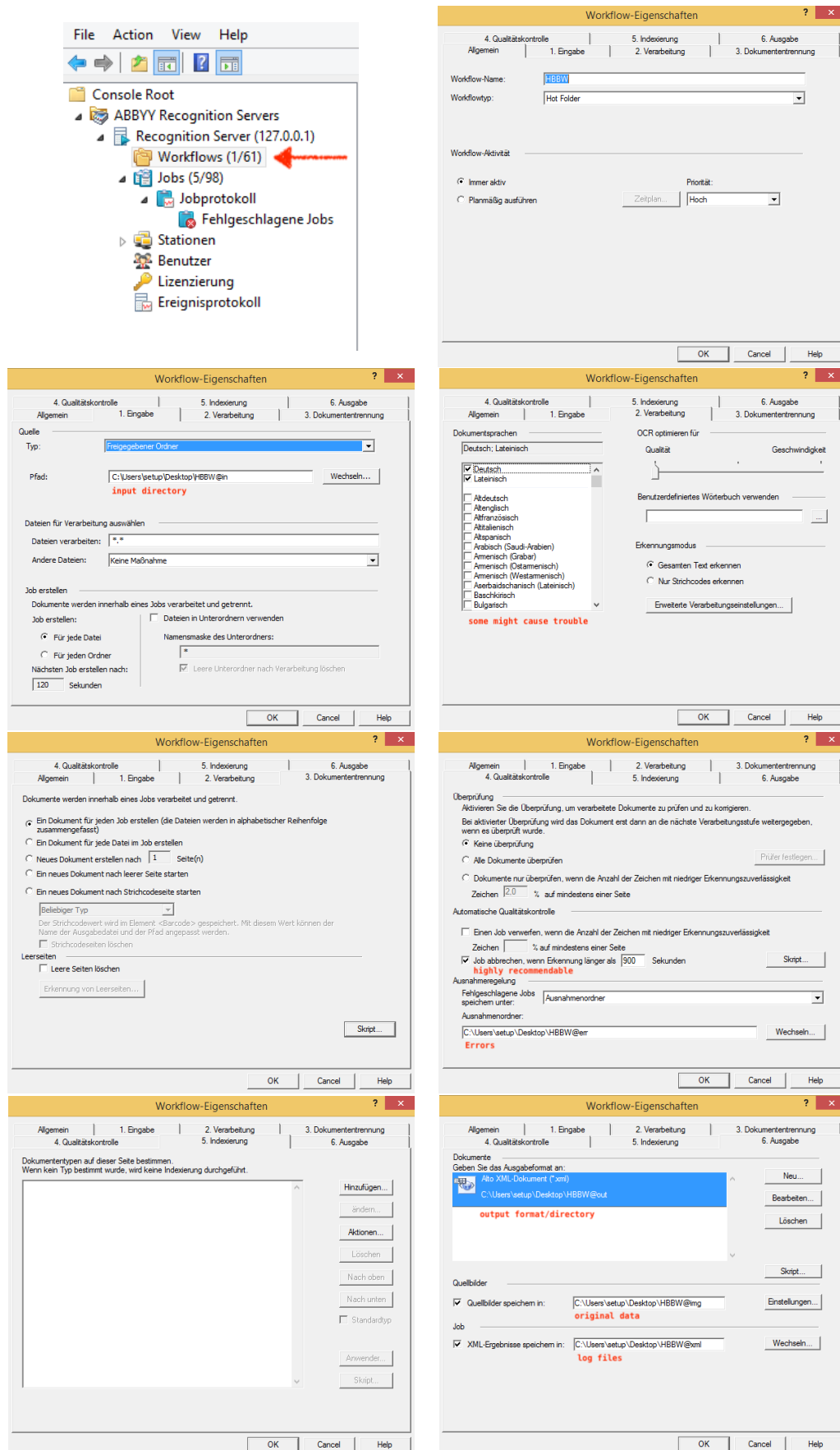
2.1 Directories (Input/Output)

FineReader requires up to 5 directories (R/W rights required):

- ▶ Input («Hot Folder»)
- ▶ Output (the OCR output in its desired format)
- ▶ Errors (corrupt input files and corresponding log files)
- ▶ Log files (optional)
- ▶ Images (optional; the input after passing the «Hot Folder»)

2.2 Configuration

The entire process can be defined in a single «Workflow»:



2.3 Execution

After completing the configuration above (by clicking on «ok»), the input directory becomes active. Every file in this directory appears to magically disappear, and (eventually) reappear again in the image directory - together with data in the other directories, as indicated by their names. If there are no configuration problems, the output follows immediately. Otherwise, the overviews «Jobs», «Job Protokoll» and «Fehlgeschlagene Jobs» (see Workflow) might reveal some information about the problem.

With the settings described, FineReader needed 3h and 49m to process our 10'012 files, thats is, 1.37s per file (80 KB) on average.

3 Kraken

The OCR-Tool «Kraken» needs to be downloaded and installed as described on the [official website](#) or on [GitHub](#). Since Kraken runs on images, our pdf files had to be transformed:

```

1 from pdf2image import convert_from_path
2 for page in convert_from_path("path/to/input_file.pdf", 600): # dpi
3     page.save("path/to/output_file.png", 'PNG')

```

Code Snipped (Python): pdf → png

```

1 import subprocess, psutil
2 def run_kraken(input_path_png, output_path):
3     cmd, url_model = "binarize segment ocr -a -m", "path/to/language_best.mlmodel"
4     cmd = ' '.join(["kraken -i", input_path_png, output_path, cmd, url_model])
5     sub_process = subprocess.Popen([cmd], shell=True)
6     p = psutil.Process(sub_process.pid)
7     try: p.wait(timeout=100)
8     except psutil.TimeoutExpired: p.kill()

```

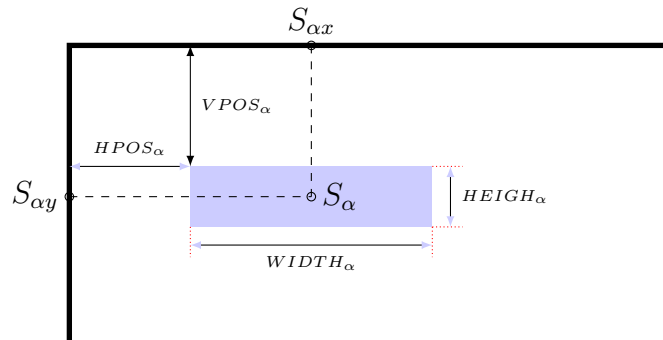
Code Snipped (Python): png → xml (alto)

4 Postprocessing

4.1 Scaling

The size of a page (length times width [px]) differs from file to file. Consequently, a reliable detection of contents based on absolute positions requires to normalize all coordinates.

- width $B_i \propto S_{i\alpha x}(\text{WIDTH}_{i\alpha x}, \text{HPOS}_{i\alpha x})$
- height $H_i \propto S_{i\alpha y}(\text{HEIGHT}_{i\alpha y}, \text{VPOS}_{i\alpha y})$



Scaling factors: $x_s(B_i) = \frac{\hat{B}}{B_i}$, $y_s(H_i) = \frac{\hat{H}}{H_i}$ ($\hat{B} :=$ norm width, $\hat{H} :=$ norm height).

4.2 Parsing

```

1 import xml.sax
2 import pandas as pd
3 from xml.sax.handler import ContentHandler
4
5 class AltoParser(ContentHandler):
6
7     """ [String, HPOS, VPOS, HEIGHT, WIDTH] → [str, x, y, baseline] """
8
9     # filter
10    NAMES = ["Datum", "Absender", "Empfaenger", "Autograph", "Kopie", "Photokopie",
11             "Standort", "Bull.", "Corr.", "Sign.", "Abschrift", "Umfang", "Sprache",
12             "Literatur", "Gedruckt", "Bemerkungen"]s
13
14    def __init__(self, scaling, inverse=False):
15        super(AltoParser, self).__init__()
16        self.data = pd.DataFrame({'Value': [], 'x': [], 'y': [], "Baseline": []})
17        self.current_baseline = 0
18        self.scaling = scaling
19        self.inverse = inverse
20
21    def startElement(self, name, attributes):
22        if name == "TextLine":
23            for a in attributes.getNames():
24                if a == "BASELINE":
25                    self.current_baseline = attributes.getValue(a)
26
27        if name == "String" and "STYLE" not in attributes.getNames():
28            d = dict()
29            d["BASELINE"] = self.current_baseline
30            for a in attributes.getNames():
31                d[a] = attributes.getValue(a)
32            if self.inverse:
33                if d["CONTENT"] not in AltoParser.NAMES:
34                    self.append(d)
35            else:
36                if d["CONTENT"] in AltoParser.NAMES:
37                    self.append(d)
38
39    def append(self, d):
40        hpos, vpos = int(d["HPOS"]), int(d["VPOS"])
41        width, height = int(d["WIDTH"]), int(d["HEIGHT"])
42        x, y = AltoParser.get_mass_point(hpos, vpos, width, height, self.scaling)
43        value = ''
44        if ("SUBS_TYPE" in d) and ("SUBS_CONTENT" in d):

```

```

45         if d["SUBS_TYPE"] == "HypPart1":
46             value = d["SUBS_CONTENT"]
47         else: value = d["CONTENT"]
48     df = pd.DataFrame({
49         'Value': [value],
50         'x': [x], 'y': [y],
51         'Baseline': [self.current_baseline]
52     })
53     self.data = pd.concat([self.data, df])
54
55     @staticmethod
56     def get_mass_point(hpos, vpos, width, height, scaling):
57         """ text center coordinates """
58         x, y = hpos + 0.5 * width, vpos + 0.5 * height
59         return int(scaling[0]*x), int(scaling[1]*y)
60
61     @staticmethod
62     def extract_data(path, scaling, inverse=False):
63         try:
64             parser = xml.sax.make_parser()
65             counter = AltoParser(scaling, inverse=inverse)
66             parser.setContentHandler(counter)
67             parser.parse(path)
68             return counter.data
69         except:
70             print("*** Warning: alto-parser failed on", path)
71             return pd.DataFrame({'Value': [], 'x': [], 'y': [], "Baseline": []})

```
