

Comparison of optical character recognition (OCR) software

by

Angelica Gabasio

Department of Computer Science
Lund University

June 2013

Master's thesis work carried out at *Sweco Position*.

Supervisors Björn Harrtell, Sweco
 Tobias Lennartsson, Sweco
Examiner Jacek Malec, Lund University

Abstract

Optical character recognition (OCR) can save a lot of time and work when information stored in paper form is going to be digitized, at least if the output from the software is accurate. If the output is very inaccurate, there will be need for a lot of post-processing, and correcting all of the errors might result in more work than manually typing in the information.

The purpose of the thesis is to run different OCR software and compare their output to the correct text to find out which software generates the most accurate result. A percentage value of the error is calculated to find out how accurate each software is. Different kinds of images will be tested to find out how accurate the software are on images of different quality.

Both commercial and open source software has been used in the thesis to find out if there is a difference.

Keywords: OCR, comparison, Tesseract, Ocrad, CuneiForm, GOCR, OCropus, TOCR, Abbyy CLI OCR, Leadtools OCR SDK, OCR API Service, Wagner-Fischer algorithm

Acknowledgements

Thanks to Björn Harrtell and Tobias Lennartsson, my supervisors at Sweco Position, for providing feedback and guidance during the thesis work.

Also thanks to everyone else at Sweco Position for being so supportive and interested in the thesis progress.

Thanks to my examiner, Jacek Malec, for feedback on this report.

Contents

1	Introduction	1
1.1	History of OCR	1
2	Theory	2
3	Method	4
3.1	Software	4
3.2	Input images	6
3.3	Comparison	6
4	Results	8
5	Discussion	12
5.1	Software	12
5.2	Results	12
5.3	Open source vs. commercial software	18
6	Conclusion	18
References		20
Appendix A: Input images		21
Appendix B: Output files		32
B.1	Tesseract	33
B.2	Ocrad	41
B.3	CuneiForm	49
B.4	GOCR	55
B.5	OCRopus	63
B.6	TOCR	70
B.7	Abbyy	75
B.8	Leadtools	80
B.9	OCR API Service	87

1 Introduction

Optical character recognition (OCR) is a way of extracting plain text from images containing text. These images can be books, scanned documents, handwritten text etc. The purpose of OCR can be to be able to search the text that has been scanned, to automatically process filled forms, or to digitize information stored in books.[1] Google books (<http://books.google.com/>) is an example of the usage of OCR scanning, where books are scanned and OCR is applied to make the text searchable.[2]

A lot of text and information is stored in paper form only, but since our surroundings is getting more and more digitized, the information would be more accessible if it was stored in digital form as well. OCR scanning can be a way to decrease the amount of manual work needed to digitize the information.

For OCR scanning to be useful it is important that the output is correct, or almost correct. The output from different software may differ a lot even if they are applied on the same image. If the output is very inaccurate there may be need for a lot of manual post-processing, and in some cases the outputted text might not be readable at all. In this thesis a comparison of different OCR tools will be made by comparing their output text to the correct text. The comparison will be done using a string comparison algorithm[4] that returns a number indicating how different two strings are. This number is used to calculate a percentage value of the output error.

1.1 History of OCR

The development of OCR scanning started in the 1950s, but there had been some similar work done earlier. In the 1960s, OCR got more widespread and the development of OCR machines increased, which lead to better accuracy of the OCR scan.

In the beginning, OCR scanning was very slow, there was a machine that read upper case alphanumeric characters at the speed of one character per minute. In the late 1950s there existed machines that could read both upper and lower case letters.[3]

In the late 1960s, there were OCR machines that could use machine learning to be able to recognize more fonts. Solutions that made it possible for the OCR machines to read characters that were imperfect, was developed at this time as well.[3]

One of the first applications of early OCR technology was the *Optophone* that was developed in 1912. It was a handheld machine that was moved over a page, and outputted different tones corresponding to the different

characters. It was developed as a help for people with vision problems, but it did not get very widespread since it required a lot from the user to learn to recognize the different tones.[3]

The first use of a commercial OCR machine in a business helped reduce the processing times of documents from one month to approximately one day, this was in 1954. At this time, OCR was also used as a help for sorting mail in post offices.[3]

One of the first OCR developers suggested the use of a standard font to help the OCR recognition in the 1950s, and in the 1960s two new fonts were introduced for this purpose. The first one, OCR-A, was designed in America, and included upper case letters, numerals and 33 other symbols. Three years later, lower case letters were added to the font as well. Europe wanted another look on the standardized font, so they designed their own standard font, OCR-B.[3]

In the beginning, the paper that was going to be OCR scanned had to be moved to scan different parts of it. The machines could do this by themselves, or either the paper or the scanner could be moved manually.[3]

2 Theory

The basic algorithm of OCR scanning consists of the following steps[1, 2]:

- Preprocessing
- Layout analysis
- Character recognition
- Output

In the *preprocessing* step the image is modified in a way that will make the OCR algorithm as successful as possible. The preprocessing usually includes:

1. Converting the image to only black and white pixels. This is called binarisation.
2. Removing noise from the binarised image.
3. Rotating the image to align the text as horizontal as possible.

In addition to these modifications there are several other less common steps that can be used for preprocessing if needed.[2] The preprocessing can be done either internally in the software or manually before the OCR scan.

The purpose of the *layout analysis* is to determine how the text should be read. This can include identifying columns or pictures in the text.[2]

After the analysis it is time for the *character recognition*, which is usually applied on one line of text at a time. The algorithm breaks the line into words, and then the words into characters that can be recognized.[2]

The algorithm uses a database of letters, numerals and symbols to match the characters. The algorithm will calculate a matching value for each character and selects the option with the best value.[1]

A way of improving the output is to check if the scanned word exists in a dictionary. If it doesn't, it is likely that some of the analyzed letters are wrong, and some of the letters may be changed so the word matches a word in the dictionary. This is done by adding a penalty value to the matching value if the analyzed character result in a word not in the dictionary. This makes it possible for another letter to get a better matching value if the word with the new letter is part of the dictionary. The letters are not always changed, sometimes a non-dictionary word gets a better matching score, even with the penalty value, than the one with the changed letters that is in the dictionary. Even if using a dictionary may help the output accuracy it is no guarantee that the correct word will be found.[1] An example of the use of a dictionary can be seen in Figure 1.

Other ways to get a better output is to use machine learning, where the OCR tool can be trained to recognize more languages or fonts, or to use the knowledge of which character combinations are more likely to appear.

The better the image is, the more accurate the result gets. If the OCR algorithm is applied on a horizontal text with very little or no noise and the letters are well separated from each other it is possible to get a very good result. But if the image is blurry, warped or noisy, or if the characters are merged, it is likely that the output gets worse than it would otherwise.[1, 2] The output accuracy will also depend on the image resolution.

OCR software is generally better at recognizing machine writing than handwriting. In machine writing the same character always looks the same, given that the same font is used. With handwriting the characters differ a lot depending on who is writing, and even if the same person has written a text, the same character will be a little different each time it is written. Machine writing often consist of more distinct characters, and since handwriting is more varied than machine writing, it is harder for the algorithm to match the characters to its database.

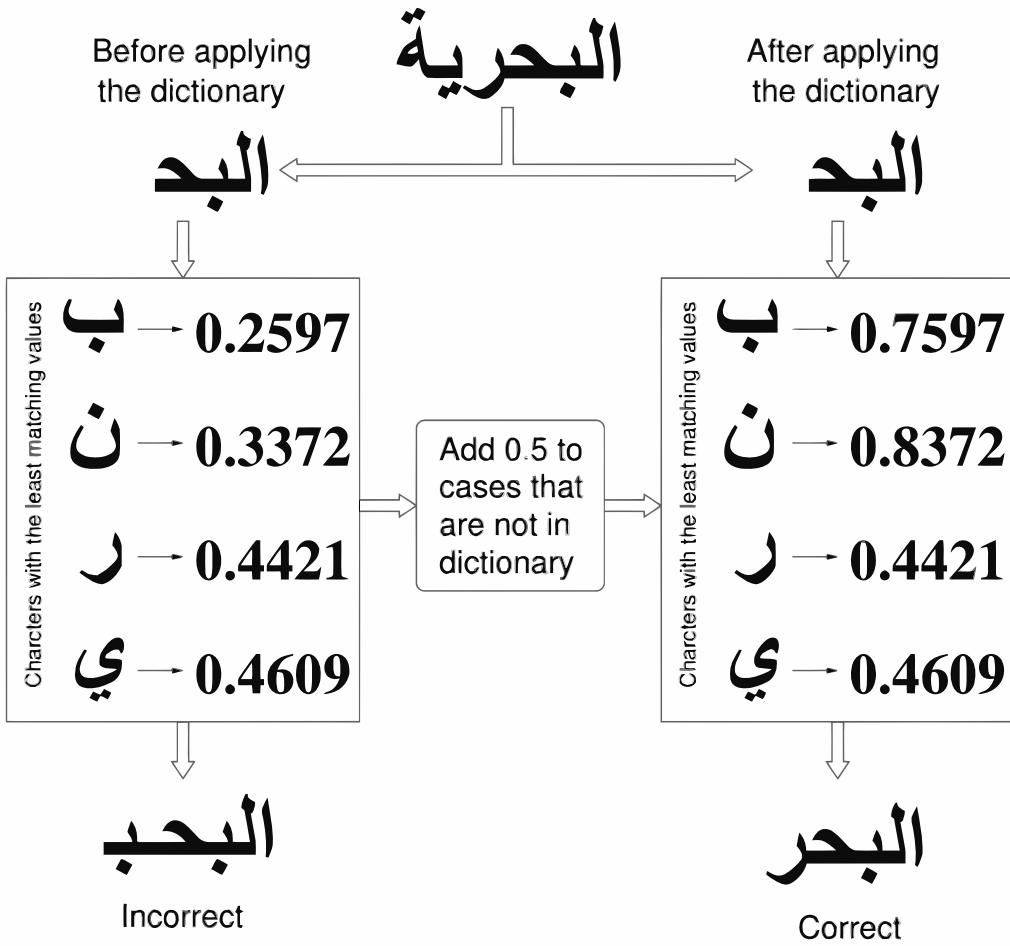


Figure 1: *The correct word is found by using a dictionary. The first letter is the one with the best matching value, but it gets a penalty since it result in a non-dictionary word. With the penalty values added, it is the third letter that has the best matching value, and results in the correct word.[1]*

3 Method

The different software used in the comparison, which images that are included, and how the comparison is done are presented here.

3.1 Software

The software that will be used in the comparison are presented below. Tesseract, Ocrad, CuneiForm, GOCR and OCropus are open source tools, TOCR,

Abbyy CLI OCR, Leadtools OCR SDK and OCR API Service are commercial tools.

The tools will not be trained for the comparison, and the settings will not be changed depending on the image, the most basic settings will be used in all of the tools to get a fair comparison. If it is possible to specify which language to use in the scans, Swedish or English will be used depending on the image.

Tesseract Tesseract is an open source OCR tool that HP started to develop in 1985. HP continued the development until 1995, and since it was released under the Apache License in 2005 the development has been sponsored by Google. Tesseract can recognize text in over 60 different languages, which have to be downloaded and added to Tesseract manually. If no language is specified for the scan, Tesseract will use English as default, Tesseract can also use multiple languages in each scan. Tesseract supports machine learning, which can be used if the language or font to be scanned doesn't exist.[5]

Ocrad Ocrad is released under the GNU GPL and is an open source OCR tool that has been developed since 2003. It is possible to define which character sets to recognize in the program, which helps narrow the character search. Ocrad can only read three different image formats so the images may need to be converted before the OCR scan.[6]

CuneiForm CuneiForm is an OCR tool that can recognize more than twenty languages and it uses a dictionary to help with the recognition. It is a tool that has been open source under the BSD license since 2008. It is built on a tool from a Russian company named Cognitive Technologies, and the documentation is in Russian. The tool only supports one image format which means that the images may need to be converted before the scan.[7]

GOCR GOCR (sometimes JOCR due to a name conflict) is an open source OCR tool released under the GNU GPL. Images may need to be converted to be used with GOCR since it can only read a few image formats.[8]

OCRopus OCROpus is an OCR tool released under the Apache License, which means it is an open source tool. OCROpus is sponsored by Google, the Mellon Foundation, and the BMBF TextGrid project, and it is mainly developed at IUPR research group. The recognition can be helped with machine learning, where OCROpus can be trained to recognize new fonts or languages.[9]

TOCR TOCR is a commercial OCR tool that has support for eleven languages, the user does not have to specify which language to use, that is done by TOCR. To help the recognition, TOCR can use character frequency/probability tables, which could help recognize misspellings. It can only read two different image formats, so there may be need for image conversion.[10]

Abbyy CLI OCR Abbyy CLI OCR is a commercial OCR tool that is based on Abbyy FineReader Engine. It supports not only horizontal text but vertical as well, which is specified when running the program. It can read 190 different languages, and offers dictionary support for some of them.[11] Abbyy offers dictionary support for both Swedish and English.

Leadtools OCR SDK Leadtools can recognize over 30 languages, and it uses spell checking and dictionary support to get a more accurate output. It is a commercial OCR tool. It is possible to get an output with the same look as the image, with the same font and layout as the image.[12]

OCR API Service This is a commercial online OCR cloud service, where the image and language is submitted through a HTTP POST request. OCR API Service has support for almost 40 different languages.[13]

3.2 Input images

Different kinds of images will be OCR scanned to see if some software is better at recognizing a specific kind of image, e.g. handwriting or noisy images. The images differ in quality, and some of them are skewed, underlined or contain pictures in the text to see how the software handles that. Most of the images contain text in Swedish, some of them are in English or a combination of Swedish and English. All of the scanned images can be seen in Appendix A.

3.3 Comparison

The comparison of the output from the OCR software to the correct text will be done using the Wagner-Fischer[4] algorithm. The algorithm measures how different two strings are by calculating the edit distance between them. The edit distance is the minimum number of edit operations that is needed for one of the strings to change into the other. The edit operations supported in the algorithm are substituting, deleting or inserting one character in one of the strings.[4] It is possible to assign different costs to the different edit

operations, but in this comparison the cost will be one for each operation. Using one as cost independent of the operation will result in a maximum edit distance equal to the length of the longer string. This corresponds to the case where every character has to be modified by one of the edit operations. If the strings are equal, no character needs to be modified and the edit distance will be zero. Pseudo code of the Wagner-Fischer algorithm with all costs equal to one can be seen in Algorithm 1, and an example of how the edit distance is calculated is shown in Figure 2.

The edit distance will be divided with the length of the longer string (the worst case edit distance), and multiplied by 100, to get a percentage value of the characters that do not match. The percentage output error that will be used for the comparison is defined by

$$\frac{d(s1, s2)}{\max\{|s1|, |s2|\}} \cdot 100$$

where d is the edit distance between the strings $s1$ and $s2$, and $|s|$ is the length of the string s .

The mean error value will also be calculated for each of the software to get an overall error value that easily can be compared.

Algorithm 1 The Wagner-Fischer algorithm with all costs equal to one[4]

Input: String $s1$, string $s2$

Output: Edit distance between $s1$ and $s2$

```

 $D[0..|s1|, 0..|s2|]$ 
for  $i = 0$  to  $|s1|$ 
     $D[i, 0] = i$ 
for  $j = 0$  to  $|s2|$ 
     $D[0, j] = j$ 
for  $i = 1$  to  $|s1|$ 
    for  $j = 1$  to  $|s2|$ 
        if  $s1[i] = s2[j]$ 
             $D[i, j] \leftarrow D[i - 1, j - 1]$ 
        else
             $d1 \leftarrow D[i - 1, j - 1] + 1$ 
             $d2 \leftarrow D[i - 1, j] + 1$ 
             $d3 \leftarrow D[i, j - 1] + 1$ 
             $D[i, j] \leftarrow \min\{d1, d2, d3\}$ 
return  $D[|s1|, |s2|]$ 

```

	a	n	g	e	l	i	c	a	
0	1	2	3	4	5	6	7	8	
g	1	1	2	2	3	4	5	6	7
a	2	1	2	3	3	4	5	6	6
b	3	2	2	3	4	4	5	6	7
a	4	3	3	3	4	5	5	6	6
s	5	4	4	4	4	5	6	6	7
i	6	5	5	5	5	5	5	6	7
o	7	6	6	6	6	6	6	6	7

Figure 2: *The edit distance between the strings is the number in the bottom right corner. The edit distance between the strings "angelica" and "gabasio" is 7.*

4 Results

The result of the OCR scans on all images is shown in Table 1. The table shows the error percentage, which means the more correct the output from the OCR is, the lower the percentage value will be. If the output text matches the correct text perfectly, the error will be 0%.

A graphical illustration of the results is shown in Figure 3. All of the output files can be seen in Appendix B. The input images are of different types, which are shown below:

Type	Image
Picture in the text	k, l
Skewed image	c, e
Handwriting	m , parts on e, f, g
Light	h
Noise	i
Stains	j

It is possible that the text in the output files is correct, but if they contain extra, or missing, whitespace that will result in an error in Algorithm 1, since it compares the strings character by character.

The mean error value, in increasing order, for each of the OCR tools is shown below:

TOCR	8.79%
Leadtools	20.06%
Abbyy	24.53%
OCR API Service	27.81%
OCRopus	31.42%
Tesseract	33.9%
CuneiForm	37.68%
Ocrad	50.25%
GOCR	64.46%

Image (shown in Appendix A)	Software				
	Tesseract	Ocrad	CuneiForm	GOCR	OCRopus
<i>a</i>	17.33%	27.45%	15.5%	51.13%	13.31%
<i>b</i>	16.78%	28.96%	23.43%	62.62%	13.86%
<i>c</i>	55.71%	40.7%	19.46%	52.04%	29.25%
<i>d</i>	31.88%	57.87%	49.01%	40.96%	13.5%
<i>e</i>	16.88%	36.2%	17.33%	77.68%	34.77%
<i>f</i>	16.4%	51.51%	17.68%	78.53%	34.34%
<i>g</i>	37.19%	75.78%	7.44%	62.86%	21.86%
<i>h</i>	45.82%	75.78%	7.44%	62.86%	21.86%
<i>i</i>	88.93%	87.72%	84.65%	92.6%	81.91%
<i>j</i>	62.79%	85.73%	8.37%	72.34%	36.36%
<i>k</i>	13.52%	28.19%	85.63%	46.47%	18.91%
<i>l</i>	2.08%	12.5%	35.5%	72.15%	37.96%
<i>m</i>	68.75%	77.59%	77.19%	91.46%	75.51%
<i>n</i>	0.57%	17.56%	78.84%	38.74%	6.49%

Table 1: Results of the OCR scans on each image from each of the open source tools.

	Software				
	TOCR	Abbyy	Leadtools	OCR API Service	
a	1.37%	8.66%	3.4%	8.82%	
b	1.6%	7.78%	3.5%	8.23%	
c	1.95%	11.51%	8.18%	7.08%	
d	1.52%	13.89%	5.01%	20.44%	
e	11.89%	16.17%	18.78%	19.3%	
f	7.83%	10.82%	23.9%	9.85%	
g	3.65%	29.07%	10.3%	53.95%	
h	3.65%	29.07%	5.80%	19.92%	
i	18.26%	100%	97.35%	100%	
j	3.65%	13.92%	28.57%	15.81%	
k	3.71%	29.95%	11.95%	19.52%	
l	3.19%	7.39%	12.27%	4.8%	
m	55.77%	63.27%	51.43%	100%	
n	4.99%	1.88%	0.44%	1.66%	

Table 1: *Continued.* Results of the OCR scans on each image from each of the commercial tools.

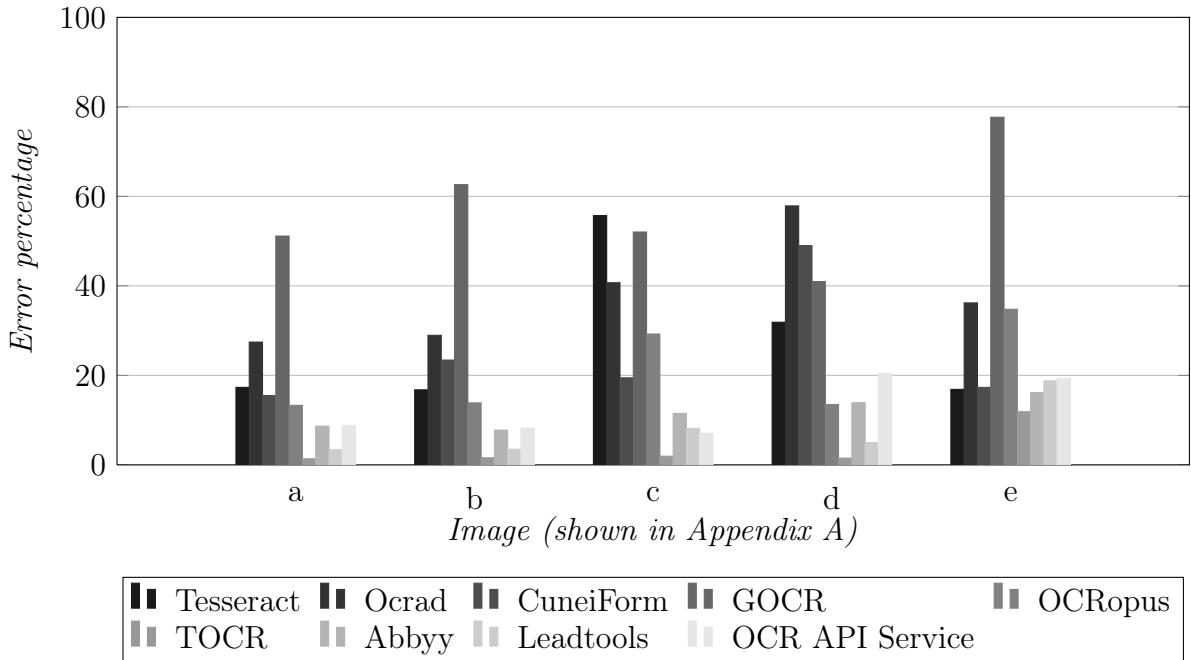


Figure 3: Illustration of the OCR scan results.

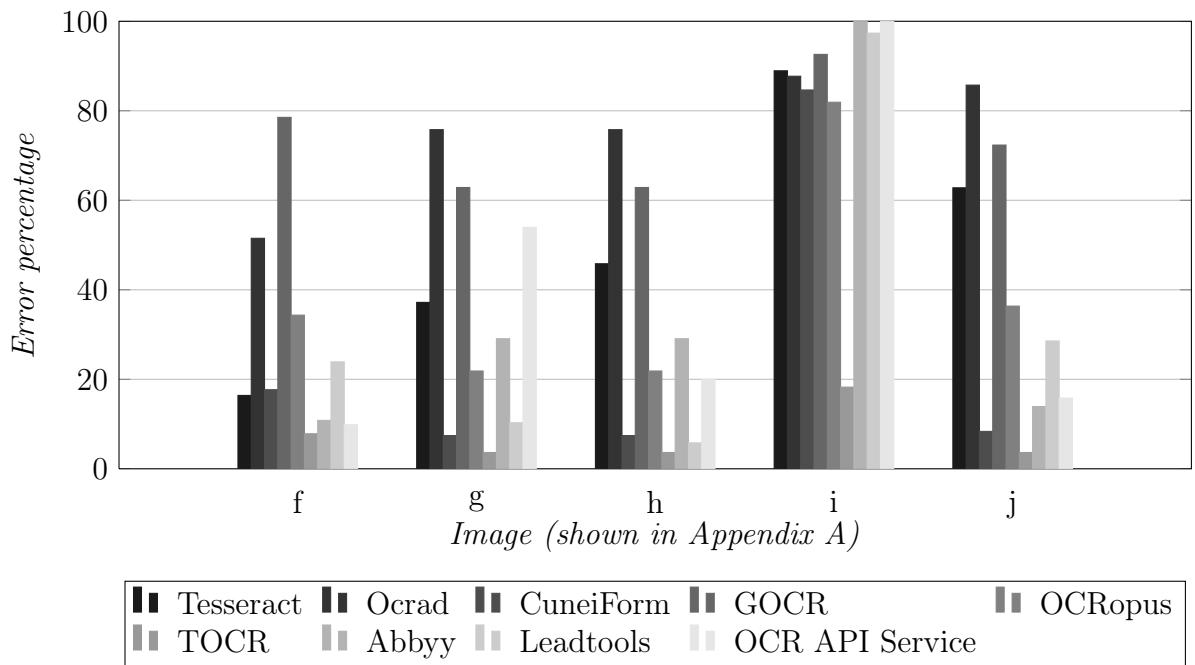


Figure 3: *Continued. Illustration of the OCR scan results.*

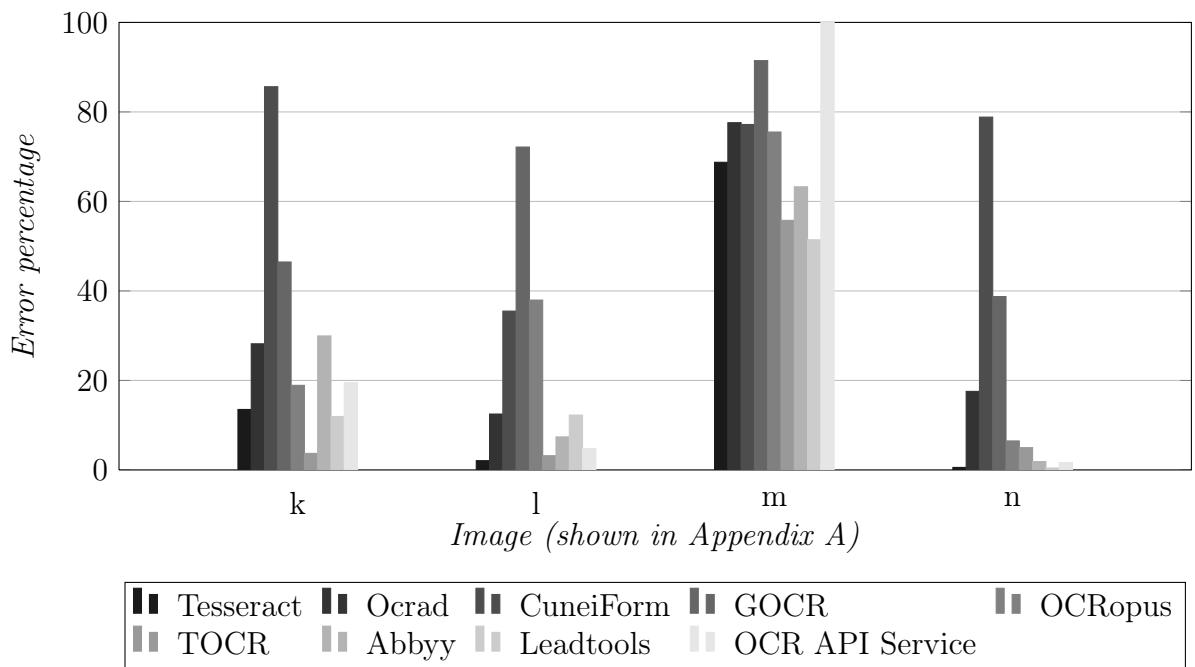


Figure 3: *Continued. Illustration of the OCR scan results.*

5 Discussion

The purpose of the thesis was to test different OCR software on different kinds of images to find out which software generates the most accurate output.

5.1 Software

Some of the commercial tools offer trial versions, limited in time or number of scans. It is difficult to limit the comparison to a certain number of scans, since some testing is needed to understand how the software works, and some of the time limited trial versions did not last long enough for the purpose. This limited the commercial tools that could be used. Emails were sent to the companies to see if they could offer a student version that could be used for the thesis. Some of them offered a version that could be used, and some just pointed at the original trial version.

It was sometimes hard to find documentation or examples on how to use the open source software. The support may also be limited when it comes to open source software, but most of them offer an email address or a forum for that.

5.2 Results

The result of each OCR scan is shown in section 4, the input images and the output files can be seen in Appendix A and B.

Image with picture Image k and image l contains pictures within the text, and the OCR tools that are best at recognizing these images are TOCR and Tesseract. On these images, TOCR has a mean error value of 3.45%, and Tesseract 7.8%.

Both TOCR and Tesseract ignore the picture in image l , and output files 82 (TOCR image l , 3.19%) and 12 (Tesseract image l , 2.08%) do not contain any extra characters where the picture is.

On image k , TOCR recognizes the numbers in the chart, but ignores the rest of it (output file 81, 3.71%). It seems like Tesseract tries to recognize the chart in image k , and output file 11 (Tesseract image k , 13.52%) contains some random characters where the chart is. TOCR also ignores the bars in the legend, Tesseract tries to recognize them and outputs some characters.

Leadtools is more accurate than Tesseract on image k with 11.95% errors compared to Tesseract with 13.52%. The output from these

tools (Leadtools: output file 109, Tesseract: output file 11) does not differ that much, Tesseract outputs more whitespace between the lines and misses most of the decimal points in the image. Since Leadtools recognizes the decimal points, I agree that the output from Leadtools is more accurate than the one from Tesseract on this image.

OCR API Service generates 19.52% errors on image k , and as seen in output file 123 (OCR API Service image k) the text is still very readable. There are some errors in the output, but most of the errors from Algorithm 1 comes from whitespace differences.

Image k contains vertical text as well, and Tesseract is the only tool that recognizes a part of it correctly as seen in output file 11.

CuneiForm and GOCR are the tools that has the highest mean error percentage on image k and l , 60.57% errors for CuneiForm and 59.31% for GOCR. Using one of those would be a bad choice if the images being scanned contain pictures.

Skewed image The skewed images are image c and e , TOCR is the tool with the most accurate result on these images, it produces a mean error percentage of 6.92%. The tool that has the highest mean error on these images is GOCR, with 64.86% errors.

Image c is a skewed version of image b , and most of the OCR tools produce a less accurate output on this image than on the original image. However, CuneiForm, GOCR and OCR API Service produce a better result on the skewed version. It is possible that it is a coincidence that the skewed image gets a lower error value. Output file 31 (CuneiForm image c , 19.46%) contains less line breaks than output file 30 (CuneiForm image b , 23.43%), and some of the lines are only recognized in the skewed case, while some other lines are only recognized in the original case. GOCR produces more whitespace between the letters in the skewed case (output file 45, 52.04%) compared to the original case (output file 44, 62.62%), which makes the text harder to read. The output texts from OCR API Service (skewed: output file 115, 7.08%, original: output file 114, 8.23%) are similar, but some of the errors are only present in one of the output texts.

Handwriting Image e and image f contains some handwritten numbers and image g contains one handwritten character. None of the OCR tools succeeds in recognizing the handwritten parts correctly on these three images. In some cases they get some of the handwritten characters right, but mostly they either ignores the handwriting or they

output what seems to be random characters.

Image m is an image with only handwriting. As seen in section 4 and in Appendix B, the output from the OCR tools is very inaccurate on this image. The software with the most accurate output on image m is Leadtools with 51.43% errors, and as seen in output file 111 the output differ a lot from the correct text.

None of the tested software claims that they can recognize handwritten text, and since all of them produce a very inaccurate output on the handwritten image, a lot of corrections need to be done in the output files. In this case it will probably be better not to use OCR scanning at all.

Noise and stains The images with added noise and stains are image h , i and j , which are all variations of image g . On the original image g , the most accurate tools are TOCR (3.65%), CuneiForm (7.44%) and Leadtools (10.3%).

Most of the OCR tools generate the same amount of errors on the lighter image h as on the original, Tesseract is the only one with more errors, 45.82% compared to 37.19% on the original image. Leadtools and OCR API Service produce better results on the lighter image than on the original. What can be seen in output file 105 (Leadtools image g , 10.3%) and 106 (Leadtools image h , 5.80%) is that the outputted text is almost the same on the two images, but the text layout differ, and the software produces some noise in the end of the text on the original image.

All of the tools are less accurate on the version with more noise, image i , and most of the error percentages differ a lot on the two images. TOCR (18.26%) is the tool with the lowest error on this image, this result is much better than the result from any of the other tools. The output from Leadtools on this image generates 97.35% errors, and output file 107 (Leadtools image i) consists only of a lot of random characters, and is not readable at all. Abbyy and OCR API Service do not recognize any of the characters in this image, and the output is blank (Abbyy: output file 93, OCR API Service: output file 121), which gives an error percentage of 100%. Even though the output on image i from Abbyy and OCR API Service results in a higher error value than Leadtools according to Algorithm 1, my opinion is that it is better to output nothing than to output lots of random characters. In this case Leadtools gets a better result according to the comparison

algorithm since some of the random characters exist in the correct text as well.

Abbyy and OCR API Service are more accurate on the stained version (image *j*) of the image than on the original (image *g*). Abbyy produces 13.92% on the stained version and 29.07% on the original image, OCR API Service produces 15.81% compared to 53.95% on the original. TOCR has the same error percentage as on the original image, while the other tools are less accurate on this image. The tools with the least amount of errors on this image are TOCR with 3.65% errors and CuneiForm with 8.37%.

Image *a* and *b* These two images are of the same document, but on image *b* there is some noise due to the paper behind showing through. On image *a*, there is no paper behind, reducing the amount of noise in the image compared to image *b*. Three of the OCR tools generate a lower error percentage on image *b*, these are Tesseract, Abbyy and OCR API Service, the rest of the tools is more accurate on image *a*. The error percentage on these two images doesn't differ that much, the tools that differ the most are CuneiForm with 15.5% errors on image *a* and 23.43% errors on image *b*, and GOCR with 51.13% errors on image *a* and 62.62% errors on image *b*.

The output from CuneiForm on image *b* (output file 30, 23.43%) misses some lines or parts of some lines in the image, which is recognized on image *a* (output file 29, 15.5%). CuneiForm recognizes "mZ" instead of "m2" in some places on image *b*, but on image *a*, the correct text "m2" is recognized. There are some other differences in the output as well, and because of the differences above, I agree that the output on image *a* is better than the one on image *b*.

GOCR is inaccurate on both of these images, and it is hard to decide if the output on image *a* (output file 43, 51.13%) really is better and easier to read than the output on image *b* (output file 44, 62.62%). Even if the result is better on image *a*, it is still more than 50% errors in the output.

Image *n* This image contains only text and minimal noise, and many of the tools produce a good and readable output on this image. Leadtools is the most accurate tool with 0.44% errors, followed by Tesseract with 0.57% errors. CuneiForm is the tool with the highest error percentage, 78.84%, and output file 42 (CuneiForm image *n*) is not readable at all and only contains random characters.

Abbyy Abbyy produces a better output on image j (13.92%) than image g (29.07%), which is a bit odd, since image j is a stained version of image g . It is possible that it is just a coincidence that the error percentage is lower, but output file 94 (Abbyy image j) does contain less noise and less whitespace than output file 91 (Abbyy image g). Another possibility to why the output is more accurate on this image may be that the software uses a different noise threshold on images with a lot of noise, and that it is able to remove more of the disturbances in the images with more noise, which results in a more accurate output.

TOCR TOCR generates the most accurate result on most of the images in the comparison. Tesseract is the only software that is more accurate on image l , where Tesseract produces an output containing 2.08% errors, and TOCR produces 3.19% errors. The output file from TOCR on this image (output file 82) contains the correct text, but it misses whitespace between the lines. Tesseract doesn't recognize the page number on this image (output file 12), but it recognizes the whitespace, which results in Tesseract getting a lower error percentage than TOCR on this image.

On image n , TOCR generates 4.99% errors, and output file 84 shows that TOCR fails to recognize the letter "d", it recognizes "cl" instead. Apart from this, the output text seems very accurate.

Since TOCR only reads two different image formats, a limitation is that images of other formats need to be converted before the scan. It might be worth converting the images and use TOCR, since it produces much more accurate results than any of the other tools in this comparison.

OCR API Service As seen in output files 121 (OCR API Service image i) and 125 (OCR API Service image m), OCR API Service does not produce any output on images i and m . Image i is the image with added noise, and image m is the handwritten image.

OCR API Service generates a more accurate result on the lighter image h (output file 120, 19.92%) and the stained image j (output file 122, 15.81%) than on the original image g (output file 119, 53.95%). It seems like the tool ignores the right side of the image in the original case (output file 119), and therefore misses some of the characters in the image. Since this issue does not appear in the lighter and in the stained version of the image, those results are more accurate.

OCRopus The most accurate open source software in this comparison is OCRopus, with a mean error of 31.42%. The tool is not the most

accurate on any of the images in Appendix A, but it is the most accurate open source tool on some of the scanned images.

The most accurate output from OCropus is 6.49% errors on image n . Output file 70 (OCropus image n) contains some recognition errors, but it is still possible to understand most of the output text.

On image i , OCropus generates 81.91% errors, which is the most inaccurate result from the tool. This is the image with added noise, and OCropus is the second most accurate software on this image, TOCR is the only tool that is more accurate with 18.26% errors. Output file 65 (OCropus image i) contains only random characters, and the output text does not match the text in the image at all.

OCropus is very slow, much slower than any of the other tools, and even though it is the most accurate open source software in the comparison, I would rather use another tool if I wanted an open source software.

Tesseract Tesseract is the second most accurate open source tool in this comparison, it produces a mean error of 33.9%. It is the most accurate tool on image l with 2.08% errors. On image n , Tesseract generates 0.57% errors, and Leadtools is the only tool with a more accurate output on this image with 0.44% errors. The text in output files 12 (Tesseract image l) and 14 (Tesseract image n) contains some recognition errors, but the text is very readable.

Tesseract is the least accurate tool on image c , which is one of the skewed images, with 55.71% errors (output file 3). The tool recognizes some parts of the image correctly, but it misses many of the lines and the output contains some noise as well.

GOCR This is the least accurate software in this comparison, with a mean error of 64.46%.

GOCR is not the least accurate software on all of the images, but it is among the least accurate tools on every image. The output files from GOCR (section B.4) contain a lot of noise, and are hard to read. The most accurate output from GOCR is 38.74% on image n , and as seen in output file 56 (GOCR image n), the text contains a lot of errors and noise and is not readable at all.

The other open source tools tested in this comparison are better choices than GOCR for OCR scanning, at least on images like the ones in this comparison.

5.3 Open source vs. commercial software

The mean error results from the comparison shows that the commercial tools are more accurate than any of the open source tools. On some images there is an open source tool that produces a lower error value, or close to the value of the commercial tools.

The most accurate open source tool is OCropus, which produces a mean error of 31.42%, compared to the least accurate commercial tool OCR API Service with 27.81% errors. Tesseract is the open source tool closest to the results of OCropus, with a mean error of 33.9%, and since the mean results does not differ that much and OCropus was very slow in the tests, Tesseract would probably be a better choice than OCropus.

The tools with the highest mean error in the comparison are Ocrad with 50.25% errors and GOCR with 64.46%. The output from these tools is very inaccurate and I think that none of the output files in Appendix B.2 (Ocrad) and B.4 (GOCR) contains an acceptable result. The most accurate result from these tools is from Ocrad on image l with 12.5% errors, and there are a lot of errors in output file 26 (Ocrad image l).

TOCR is the software with the most accurate output (mean error 8.79%), the tool that is closest to this result is Leadtools with a mean error of 20.06%. The mean error percentages from these tools are not that close, and TOCR generates a more accurate result than Leadtools on most of the images in Appendix A, Leadtools is more accurate than TOCR on two images, image m and n . The error percentages from TOCR and Leadtools on image m (TOCR: 55.77%, Leadtools: 51.43%) does not differ that much, and the output from both of the tools is inaccurate and very hard to read (TOCR: output file 83, Leadtools: output file 111). On image n , TOCR generates 4.99% errors and Leadtools 0.44%, this is the image where TOCR fails to recognize the "d's in the text (output file 84). Since Leadtools (output file 112) does not have this problem, that output is more accurate than the one from TOCR.

6 Conclusion

As seen in section 4, the OCR tool with the lowest mean error value on all of the images in Appendix A, is the commercial software TOCR with 8.79% errors. The software closest to the results of TOCR is the commercial tool Leadtools, with 20.06% errors. The most accurate open source software is OCropus with a mean error value of 31.42%, followed by Tesseract with 33.9% errors.

The four commercial OCR tools are the ones that have the lowest mean error values, this shows that it probably would be a good idea to invest in a commercial OCR software, at least if the images being scanned are of different quality, as the ones in this thesis. The results from TOCR and the other tools differ a lot and the better choice would be TOCR.

If a commercial tool is unwanted, the mean error shows that OCropus is the most accurate open source tool, but since it is very slow, Tesseract, which is the second most accurate open source tool, would be a better choice, at least if the scanning time is an issue. The mean error does not differ that much between these tools, OCropus generates a mean error of 31.42% and Tesseract 33.9%.

Ocrad and GOCR are the software that has the highest mean error values (Ocrad: 50.25%, GOCR: 64.46%) in this comparison, using one of these tools would be a bad choice if the images that is going to be OCR scanned are similar to the ones in Appendix A.

References

- [1] Inad Aljarrah, Osama Al-Khaleel, Khaldoon Mhaidat, Mu'ath Alrefai, Abdullah Alzu'bi, Mohammad Rabab'ah, *Automated System for Arabic Optical Character Recognition with Lookup Dictionary*, Journal of Emerging Technologies in Web Intelligence, Nov 2012, Vol. 4 Issue 4, pp. 362-370
- [2] Tobias Blanke, Michael Bryant, Mark Hedges, *Open source optical character recognition for historical research*, Journal of Documentation, Vol. 68 Iss: 5, pp. 659-683
- [3] Herbert F. Schantz, *The history of OCR: Optical character recognition*, Recognition Technologies Users Association, 1982
- [4] Robert A. Wagner, Michael J. Fischer, *The String-to-String Correction Problem*, Journal of the Association for Computing Machinery, Vol. 21, No. 1, January 1974, pp. 168-173

Software used

- [5] Tesseract-ocr, <http://code.google.com/p/tesseract-ocr/>
- [6] Ocrad, <http://www.gnu.org/software/ocrad/>
- [7] CuneiForm, <http://cognitiveforms.com/ru.html#1189-CuneiForm>
- [8] GOCR, <http://jocr.sourceforge.net/>
- [9] OCropus, <https://code.google.com/p/ocropus/>
- [10] TOCR, <http://www.transym.com/tocr-the-integrators-choice.htm/>
- [11] Abbyy CLI OCR, <http://www.ocr4linux.com/>
- [12] Leadtools OCR SDK, <http://www.leadtools.com/sdk/ocr/>
- [13] OCR API Service, <http://ocrapiservice.com/>

Appendix A: Input images

Forts
Fastigheten Bilen 4

6.

Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
Del om 836 m^2 av fastigheten Bilen 4 utgör nybildning av samfälligheten Stapelbädden S:4 (se överenskommelse i avyttringsakt ang del om ca $3\ 020 \text{ m}^2$ av fastigheten Bilen 4)
Areal: $339\ 546 \text{ m}^2$

Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737 m^2 av fastigheten Bilen 4 utgör Galeasen 1.
Del om $2\ 229 \text{ m}^2$ av fastigheten Bilen 4 utgör Galeasen 2.
Areal: $336\ 580 \text{ m}^2$

V G V!

Forts från fastigheten
Bilen 4

Sid 7

Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 m^2 av samfälligheten Flaggskepparen S:1 har tillagts fastigheten Bilen 4.
Del om 9 m^2 av fastigheten Bilen 4 har tillagts samfälligheten Flaggskepparen s:1.
Areal: $332\ 894 \text{ m}^2$

Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om $96\ 478 \text{ m}^2$ av fastigheten Bilen 4 utgör fastigheten Bilen 12.
Areal: $236\ 416 \text{ m}^2$

V G V!

Image a

Forts
Fastigheten Bilen 4

6.

Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
Del om 836 m^2 av fastigheten Bilen 4 utgör nybildning av
samfälligheten Stapelbädden S:4 (se överenskommelse i avyttrings-
akt ang del om ca $3\ 020 \text{ m}^2$ av fastigheten Bilen 4)
Areal: $339\ 546 \text{ m}^2$

Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737 m^2 av fastigheten Bilen 4 utgör Galeasen 1.
Del om $2\ 229 \text{ m}^2$ av fastigheten Bilen 4 utgör Galeasen 2.
Areal: $336\ 580 \text{ m}^2$

V G V!

Forts från fastigheten
Bilen 4

Sid 7

Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 m^2 av samfälligheten Flaggskepparen S:1 har tillagts
fastigheten Bilen 4.
Del om 9 m^2 av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s:1.
Areal: $332\ 894 \text{ m}^2$

Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om $96\ 478 \text{ m}^2$ av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: $236\ 416 \text{ m}^2$

V G V!

Image b: Almost the same as image a, but some text from another page is visible through the paper.

Forts
Fastigheten Bilen 4

6.

Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
Del om 836 m^2 av fastigheten Bilen 4 utgör nybildning av
samfälligheten Stapelbädden S:4 (se överenskommelse i avyttrings-
akt ang del om ca $3\ 020 \text{ m}^2$ av fastigheten Bilen 4)
Areal: $339\ 546 \text{ m}^2$

Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737 m^2 av fastigheten Bilen 4 utgör Galeasen 1.
Del om $2\ 229 \text{ m}^2$ av fastigheten Bilen 4 utgör Galeasen 2.
Areal: $336\ 580 \text{ m}^2$

V G V!

Forts från fastigheten
Bilen 4

Sid 7

Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 m^2 av samfälligheten Flaggskepparen S:1 har tillagts
fastigheten Bilen 4.
Del om 9 m^2 av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s:1.
Areal: $332\ 894 \text{ m}^2$

Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om $96\ 478 \text{ m}^2$ av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: $236\ 416 \text{ m}^2$

V G V!

Image c: A skewed version of image b.

Enligt avtal A 515 och F 11 har del om ca 760 m² av fastigheten
2009
Bilen 4 avyttrats. Se akt ang delar om ca 2 060 m² resp 11 m²
av fastigheten Bilen 7 i lagf skåp. Avst reg 16/3-2011

Enligt avtal A 34/2010 har del om ca 2 400 m² av fastigheten
Bilen 4 avyttrats. Avst reg 15/7-10

Avst reg 15/7-10: Se avyttringsakt ang del om ca 2 400 m² av
Bilen 4.
Del om 2 428 m² av Bilen 4 utgör fastigheten Klyvaren 1. A 34/201
Del om 1 312 m² av Bilen 4 utgör fastigheten Klyvaren 2.
Areal: 332 840 m²

Forts sid 7

Avst reg 25/1-2011: Se akt ang Bilen 4 i lagf skåp
Del om 1 017 m² av Bilen 4 utgör Klippen 1.
Del om 813 m² av Bilen 4 utgör Klippen 2.
Del om 631 m² av Bilen 4 utgör Klippen 3.
Del om 1 078 m² av Bilen 4 utgör Klippen 4.
Del om 2 255 m² av Bilen 4 utgör Kosterbåten 1.
Del om 1 627 m² av Bilen 4 utgör Kosterbåten 2.
Del om 1 784 m² av Bilen 4 utgör Koggen 1.
Del om 1 666 m² av Bilen 4 utgör Koggen 2.
Areal: 225 545 m²

Avst reg 16/3-2011: Se akt ang delar om ca 2060 m² resp 11 m² av
fastigheten Bilen 7 i lagf skåp

Del om 759 m² av Bilen 4 utgör Bilen 13. A 515 och F 11
Areal: 224 786 m²

2009

Image d

Fastigheten Oxie 1:5

Frgl reg 19/4-79: se akt ang fastigh Oxie 1:1 m fl
Del om 10 801,7 m² av samfälligheten Oxie s:2 har tillagts
fastigheten Oxie 1:5 (F 29/1979)
Å 27/1979

2234 580,5 m²

Enl. avtal A61, F60 1979 har del av fastigheten Oxie 1:5 avyttrats,
se akt ang. ca 78.000 m² av fastigheten Oxie 21;2. Delen har
vid frgl. reg. 15/1-81 visat sig innehålla 7.287,9 m².

Areal: 2 227 292,6 m²
Enl. avtal A 20/1981 har del om 32.615,9 m² av Oxie 1:5, avs.
att bilda fastigheten Fornlämningen 1 avyttrats. (avsl reg 17/1-81)

Avst reg 17/3-88: se akt ang Oxie 1:1 m fl
Del om 14 225 m² ay Oxie 1:5 utgör Fornlämningen 2
Areal: 2 180 452 m²

Image e

657,654

Fastigheten Oxie 1:5

Bildad genom sammanläggning av fastigheterna Oxie 1:1,
Oxie 17:2, Oxie 19:1, Oxie 20:8, Oxie 25:9, Oxie 28:1 och
Oxie 29:1.

Areal efter sammanläggning: 2.231.222 m²

Särskilt namn enl. reg. beslut den 15/4-92: Fredriksberg
Frgl. reg. 17/10-75 : se akt ang. stg 2440 m.fl. i Fosie
Delar om 85,0 m² av samfällda området adj, 1.956,0 m² av sam-
fällda området Oxie s:1, 33.507,2 m² av stadsägan 2834 och
22.122,4 m² av stadsägan 2467, 2468 har tillagts fastigheten
Oxie 1:5.

Image f

Forts från
Fastigheten Oxie 1:5

4.

Del om 1 261 m² av fastigheten Oxie 1:5 utgör samfälligheten

Gånggriften S:1.

Areal: 2 041 038 m²

Enligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten

Oxie 1:5 avyttrats.

1

Image g

Forts från
Fastigheten Oxie 1:5

4.

Del om 1 261 m² av fastigheten Oxie 1:5 utgör samfälligheten

Gånggriften S:1.

Areal: 2 041 038 m²

Enligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten

Oxie 1:5 avyttrats.

1

Image h: A lighter version of image g.

Forts från
Fastigheten Oxie 1:5 4.

Del om $1\ 261\ m^2$ av fastigheten Oxie 1:5 utgör samfälligheten

Gånggriften S:1.

Areal: $2\ 041\ 038\ m^2$

Enligt avtal A 329/2006 har del om ca $23\ 300\ m^2$ av fastigheten
Oxie 1:5 avyttrats. 1

Image i: A noisier version of image g.

Forts från
Fastigheten Oxie 1:5 4.

Del om $1\ 261\ m^2$ av fastigheten Oxie 1:5 utgör samfälligheten

Gånggriften S:1.

Areal: $2\ 041\ 038\ m^2$

Enligt avtal A 329/2006 har del om ca $23\ 300\ m^2$ av fastigheten
Oxie 1:5 avyttrats. 1

Image j: A stained version of image g.

Software	Image (shown in Appendix A)				
	a	b	c	d	e
Tesseract	16.97%	35.59%	16.72%	36.78%	16.4%
Ocrad	28.96%	57.93%	36.2%	75.78%	51.51%
OpenOCR	23.43%	49.01%	17.33%	7.442%	17.68%
GOCR	62.62%	41.04%	77.68%	62.86%	78.53%
OCRopus	32.87%	15.79%	37.64%	41.52%	39.38%
TOCR	4.74%	4.413%	13.22%	7.207%	11.26%
Abbyy	9.513%	29.16%	24.56%	27.03%	15.31%

Table 1: Results of the OCR scans.

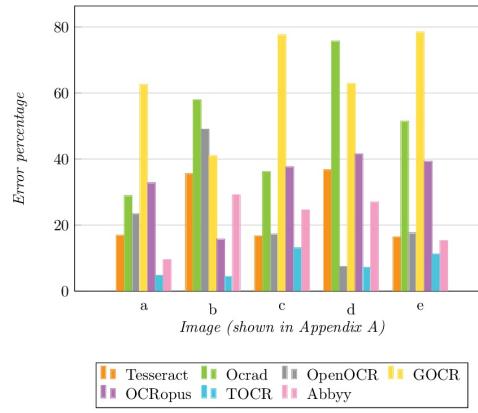


Figure 2: Illustration of the scan results.

The mean error value in increasing order for each of the OCR tools is shown below:

TOCR	8.169%
Abbyy	21.11%
OpenOCR	22.98%
Tesseract	24.49%
OCRopus	33.44%
Ocrad	50.07%
GOCR	64.55%

Testing OCR tools

Angelica Gabasio

1 English

This is a page with a picture to test the OCR tools.
The page contains text in English and Swedish.



Figure 1: This is a cat.

2 Swedish

Och så lite svenska tecken. Båten åker i sjön.
Det är blött.

1

Image l

Optical Character recognition

ANGELICA GABASIO

2013

Image m

The algorithm uses a database of letters and numerals to match the characters. The algorithm will calculate a matching value for each character and selects the option with the best value.[1]

A way of improving the output is to check if the scanned word exists in a dictionary. If it doesn't, it is likely that some of the analyzed letters are wrong, and some of the letters may be changed so the word matches a word in the dictionary. This is done by adding a penalty value to the matching value if the analyzed character result in a word not in the dictionary. This makes it possible for another letter to get a better matching value if the word with the new letter is part of the dictionary. The letters are not always changed, sometimes a non-dictionary word gets a better matching score, even with the penalty value, than the one with the changed letters that is in the dictionary. Even if using a dictionary may help the output accuracy it is no guarantee that the correct word will be found.[1] An example of the use of a dictionary can be seen in Figure 1.

Other ways to get a better output is to use machine learning, where the OCR tool can be trained to recognize more languages or fonts, or to use the knowledge of which character combinations are more likely to appear.

The better the image is, the more accurate the result gets. If the OCR algorithm is applied on a horizontal text with very little or no noise and the letters are well separated from each other it is possible to get a very good result. But if the image is blurry, warped or noisy, or if the characters are merged, it is likely that the output gets worse than it would otherwise.[1, 2] The output accuracy will also depend on the image resolution.

OCR software is generally better at recognizing machine writing than handwriting. In machine writing the same character always looks the same, given that the same font is used. With handwriting the characters differ a lot depending on who is writing, and even if the same person has written a text, the same character will be a little different each time it is written. Machine writing often consist of more distinct characters, and since handwriting is more varied than machine writing, it is harder for the algorithm to match the characters to its database.

Appendix B: Output files

The output from the OCR tools on the images in Appendix A is shown here.

Some of the output files contain rows that are too wide to fit on the page, these files have been forced to break those lines. The original output files, with no extra line breaks, are the ones used in the comparison.

B.1 Tesseract

Forts a- 6a

FastigheteáfBilen 4 *\

Frgi reg 19/8-09: Se akt ang Bilen 4 i 1agf skåp

Del om 836 m² av fastigheten Bilen 4 utgör nybiidning av

samfälligheten Stapeibädden S:4 (se överenskommelse i avyttrings- .

2

Ä akt ang de] om ca 3 020 m av fastigheten Bilen 4)

ä Areaiz 339 546 m²

Avst reg 0 /9-09:; se akt ang Bilen 4 i 1agf skåp
5 Del om 737 m² av fastigheten Bilen 4 utgör Gaieasen 1;

De] om 2 229 m² av fastigheten Bilen 4 utgör Galeasen g-
Areal: 336 580 m²,

V G V!

i

I

1

1

Forts från fastigheten

2 Bilen 41 w m" m VVm m ww
V Frgl reg 30/7-10: Se akt ang Bilen 4 i 1agf skåp

De] om 63 mz av samfälligheten Flaggskiepparen 5:1 har tiliagts
fastigheten Bilen 4.

2 av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s:1.

Areal: 332 894 m²

Del om 9 m

Avst reg 23/11-2010: Se akt ang Bilen 4 i 1agf skåp

D61 Om 95 478 m² av fastigheten Bilen 4 utgör fastigheten
Bilen 12.

Areal: 236 416 m²

V G V!

Output file 1: *Tesseract – Image a*

å
va,

r a' 'my' å\\
r

Forts A, 2 x
n. Bilen 4

, \
Fastighete

f' .\
l

Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
Del om 836 m² av fastigheten Bilen 4 utgör nybildning av
samfälligheten Stapelbädden S24 (se överenskommelse i avyttrings- .
akt ang del om ca 3 020 m²
Areal: 339 546 m²
av fastigheten Bilen 4)

Avst reg 0 /9-09 i; Se akt ang Bilen 4 i lagf skåp
Del om 737 m² av fastigheten Bilen 4 utgör Galeasen 1;

Del om 2 229 m²
Areal: 336 580 m²

av fastigheten Bilen 4 utgör Galeasen g-
/
V G V!

Forts från fastigheten Sid 7
Bilen 4- __mⁿ - -

J

Frgl reg 3073-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 mz av samfälligheten Flaggskepparen 5:1 har tillagts
fastigheten Bilen 4.

Del om 9 m² av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s:1.

Areal: 332 894 m²

Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
D61 Om 95 478 M² av fastigheten Bilen 4 utgör fastigheten

Bilen 12.
Areal: 236 416 m²

V G V!

Output file 2: *Tesseract – Image b*

Forts , , '% \;
Fastigheten/Bilen 4 1

6.
1_____g_____m-

Frgl reg 19/8-09: Se akt ang Bilen 4 i l

Del om 836 m2

V G V'

Forts från fastigheten \$'d 7

Å'l1i MM . .M .
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf sk
Del om 63 m2 av samfälli

fastigheten Bilen 4.

Del om 9 m2 av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s:1.

Areal: 332 894 m2

I Avst reg 23/11-2010:

Del Om 95 478 M2 av fastigheten Bilen 4 utgör fastigheten
Bilen 12.

Areal: 236 416 m?

Output file 3: *Tesseract – Image c*

. ?Eli F_ Vta] A 515 Ch F 11 har de? om ca 760 m2 av fastighgten
- .. ?009 W .
- 3116" êgvxtratsaz --S9. akatansadeiaráom ca 2 060 m2 resp li mf
av fa5ti9" e t e 3ÅJ?U-7 _ T lâ9f. âkâp-wÅY\$t_f69.16 /3'2 Qll

gu . Enligt avtal A 34/2010 ihar de? om ca 05400 m2 av fasttg eten
d, Bil?" 4 .avytrats. Aysj; reg 15/7710 ____m_____.
:gÅJstHrEg55Alá77-1qgg5 Se 5 \$ ttr '1"ägsá L é ang 221 omu ca? 460 in? av
0 Bilen_4, " , ä_

I 0e1_om z_42e mz ;Q äilep 4 utgår fastjghetéa klyváké 1. Ä 34/261

De1.om.1.31 z m2_av Bilen_4.utgörwfastighetgn Klyvaren 2:
_ Area]: _332g840 m2_

vn-M - - - u - - HM --N--4h__4wmwlm__4t_____m_- Eorts sidM7

.I Del om 1 017 m2 av Bilen 4 utgör Klipern 1.

..J Del om 813 m² av Bilen 4 utgör Klippern 2.

{ Del om 631 m² av Bilen 4 utgör Klippern 3.
I Del om 1 078 m² av Bilen 4 utgör Klippern 4.

I Del om 2 255 m² av Bilen 4 utgör Kosterbåten 1.

2
f Del Om 1 627 nä-av Bilen 4 utgör Josterbåten 2.
rDel om 1 784 m av Bilen 4 utgör Koggen 1.

.HI D91 m 1 555 m² av Bilen 4 utgör Koggen 2.
Areal: 225 545 m²

Avst reg 16/3-2011: Se akt ang delar om ca 2060 m² resp 11 m² av
åfastigheten Bilen 7 i 1agf skåp
| 0%: om 759 m² av Bilen 4 utgör Bilen 13. A 515 och F 11

' Areal: 224 786 m² 2 0 0 9

4...

Output file 4: *Tesseract – Image d*

Fastigheten Oxie 1:5

Frgl reg 1214-79: se akt ang fastigh Oxie 1:1 m fl

Del om 10 801,7 Å av samfälligheten Oxie s:2 har tillagts
fastigheten Oxie 1:5 (E å3%13;g)
Z.å3V.5?O,5'07

A61 F60
Enl- aVtal "J", har del av fastigheten Oxie 1:5 avyttrats,e

1979
se akt ang. ca 78.000 Å av fastigheten Oxie 21;2. Delen har
vid frgl. reg. 15/1-81 visat sig innehålla 7.287,9 E.

Areal: 2 227 292,6 f
Enl. avtal A 20/1981 har del om 52.615,9 Å av Oxie 1:5, avs.

(MJ M7 ?D 'KID

att bilda fastigheten Fornlämningen 1 avyttrats.

Avst 19% 17/3-88: se akt ang Oxie 1:1 m fl

Del om lü 225 m² å Oxie 1:5 utgör Fornlämningen 2

Axsai: 2 180 u52 m

Output file 5: *Tesseract – Image e*

Q57, Q V

Fastigheten Oxie 1:5

Bildad genom sammanläggning av fastigheterna Oxie 1:1,

Oxie 17:2, Oxie 19:1, Oxie 20:8, Oxie 25:9, Oxie 28:1 Och
Oxie 2931.

;Areal' effer sam@anlä" nin : 2.231.222 å

Sä skiltinamn -eni; reg beslut den 15/4f921. EL dIÅk D LQ1

Frgl, reg.17{10 _Z : se akt ang. stg 2440 m. fl. i Fosie

Delar om 85,0 J av samfällda området adj, 1.956,0 f av sam-
fällde området Oxie s:1. 35.507,2 ä av stadsägan.2834 och

22.12?,4 Å av seadsägan 9å67,2468 har tilla ts fastigheten

Output file 6: *Tesseract – Image f*

W Forts från

Fastigheten Oxie 1:5 4.

-D61 9m 1?261 iäwav faSti9heten ,0 ieh1 =5 utgör mf ä11i9h ten_-
_Eänifi fteg! 5:1-,
mAreal: 2 041 038 m2

ää1igtZävtal Ä 329[gbÖê:Har;gê1ZgåWcg_23 300 m? gyqfa tigheten
"Oxie 1:5 avyttra; ;, 4

Output file 7: *Tesseract – Image g*

W Forts från

Fastigheten Oxie 1:5 4.

_De1 pm 11261 m "ay fastjghe eg_O je 1:5 utgör mf ä11jghçtegn<
_Qånggtiftegy 5:1. - -{ ,i ,
mAreal: 2 041 038 m2

wiå kligtw åvtaJjÄ ;zái/:z öhdêiwrrarf _åêl Garcia: 2; swmzwçfvqaiastigwsn-
"Oxie 1:5 avyttra; ;, 4

Output file 8: *Tesseract – Image h*

?mig ?rää V V A
?väääääääüäéiâ üxêü å.

wvxer - w - snw , muwwmwmmww mw mm xw wnu m wwamwwm a . wwwmww rwá
vwuywwmsm&w m c aiw äwwömähwbwemwnwwuamwwwww

\<
, w ? á I,
%, v.& a å .;.*: ,; 2

Output file 9: *Tesseract – Image i*

```
di* 5:7  
Jb -W  
1 Forts fr n < ::$.  
Fastigheten Oxie 1:5 4.
```

```
De1 om 1 261 ;% ay fastigheten Oxie 1:5 utg rwggf llighetenz  
Q ngggiften Szltrz  
_Areaiz 2 041 038 _mz
```

Output file 10: *Tesseract – Image j*

Image (shown in Appendix A)

```
a b C d e  
Tesseract 1697% 3559% 1672% 3678% 164%  
Ocrad 2896% 5793% 362% 7578% 51151%  
  
OpenOCR 2343% 4901% 1733% 7442% 17168%  
GOCR 6262% 4104% 7768% 6286% 78.53%  
OCROpus 3287% 1579% 3764% 41152% 39.38%  
TOCR 474% 4413% 1322% 7207% 1126%  
Abbyy 9513% 29.16% 2456% 2703% 1531%
```

Software

Table 1: Results of the OCR scans.

```
so  
0 60  
,3
```

U

```
3 40  
E  
La
```

20

0

a b c d e

Image (shown in Appendix A)

```
Iii Tesseract Ii Ocrad lili OpenOCR Iii GOCR  
OCROpus In TOCR M Abbyy
```

Figure 2: Illustration of the scan results.

The mean error value in increasing order for each of the OCR tools is shown below:

```
TO CR 81169%  
Abbyy 21.11%  
OpenOCR 22198%  
Tesseract 24149%  
OCROpus 33144%  
Ocrad 5007%  
G 0 CR 6455%
```

Output file 11: *Tesseract – Image k*

Testing OCR tools

Angelica Gabasio

1 English

This is a page with a picture to test the OCR tools.
The page contains text in English and Swedish.

Figure 1: This is a cat.

2 Swedish

Och så lite svenska tecken. Båten åker i sjön.
Det är blött.

Output file 12: *Tesseract – Image l*

```
OP\:\i(.Q\ C\m\rm:5:e.r racoqn\ ion  
ANe1:\_\c/+\\ (';{A%/\5\  
2075
```

Output file 13: *Tesseract – Image m*

The algorithm uses a database of letters and numerals to match the characters. The algorithm will calculate a matching value for each character and selects the option with the best value.[1]

A way of improving the output is to check if the scanned word exists in a dictionary. If it doesn't, it is likely that some of the analyzed letters are wrong, and some of the letters may be changed so the word matches a word in the dictionary. This is done by adding a penalty value to the matching value if the analyzed character result in a word not in the dictionary. This makes it possible for another letter to get a better matching value if the word with the new letter is part of the dictionary. The letters are not always changed, sometimes a non dictionary word gets a better matching score, even with the penalty value, than the one with the changed letters that is in the dictionary. Even if using a dictionary may help the output accuracy it is no guarantee that the correct word will be found.[1] An example of the use of a dictionary can be seen in Figure 1.

Other ways to get a better output is to use machine learning, where the OCR tool can be trained to recognize more languages or fonts, or to use the knowledge of which character combinations are more likely to appear.

The better the image is, the more accurate the result gets. If the OCR algorithm is applied on a horizontal text with very little or no noise and the letters are well separated from each other it is possible to get a very good result. But if the image is blurry, warped or noisy, or if the characters are

merged, it is likely that the output gets worse than it would otherwise.[1 ,
2]
The output accuracy will also depend on the image resolution.

OCR software is generally better at recognizing machine writing than handwriting. In machine writing the same character always looks the same, given that the same font is used. With handwriting the characters differ a lot depending on who is writing, and even if the same person has written a text, the same character will be a little different each time it is written.
Machine writing often consist of more distinct characters, and since handwriting is more varied than machine writing, it is harder for the algorithm to match the characters to its database.

Output file 14: *Tesseract – Image n*

B.2 Ocrad

| ' Fartg ,__ .-' . i \ S -
rastighe_e_ 8ilen_ 4 ;,-\ -
. : -+ . - - - . - - - - - -
- : -Se akt ang 8ilen 4 i lagf skäp
- Del om 836 m' av rastigheten Bilen 4 utgör nybildning av
- samrä17igheten Stapelbädden s:4 (se överenskormP7se i avytrrings -
- akt ang del om ca 3 ozo m' av fastigheten Bilen 4)
Areal: 339 5q6 m'
i -

_1 ___, 'Se akt ang Bi7en 4 i lagf skäp
De7 om 737 'mZ...aY rastigheten Bilen 4 utgör Galeasen l...
Del om Z zz9 m' av rastigheten Bilen 4 utgör Caleasen z_
: 336 580 mZ "
|
| y G v!
|
|
|
| Forts frán fasti héten sid _
| Bilen 4
| k_ à g Bi | é _ - - . fä - - - - - - - -
De7 om 63 m' av samfälligheten Flaggsképparen s:l har tillagts
fastigheten Bilen 4.
Del om 9 m' av rastigheten Bilen 4 har ti7lagts samrä)7igheten
Flaggskepparen s:l.
Areal: 33z 894 m'

: Se aktang Bilen 4 i lagf skäp
Del om 96 478 m' av fastigheten Bilen 4 utgör fastigheten
8ile.n IZ.
Areal : 236 416 m2
y G v!

Output file 15: *Ocrad – Image a*

_ . fartg , , / - ' . - - ; , --- , 6 ,
_ astiahe _ e _ Bilen 4 , - - ; | ---

- : Se aktang Bilen4 i lagr skäp
Del om 836 m' av rastigheten Bilen 4 utgör nybi7dning av
, samrä7ligheten Stapelbädden s:4 (se överenskomnelse i avyttrings-
akt ang de) om ca 3 ozo m' av fastigheten Bi7en 4)
i Areal : 339 546 m'
|
_____; 'Se akt ang Bilen 4 i lagf skäp
Del om 73 mZ.ay rastigheten Bilen 4 utgör Galeasen 1 ..
Del om Z zz9 m' av rastigheten Bilen 4 utgör GaFeasen z_
Areal : 336 580 mZ "
i -
|
v c v!
| .
|

i Forts frgn fasti héten sid _
Bilen 4
- : Se ak äng 8ile _ - . - íä r -- - p - - -
Del om 63 m' av samfäl)igheten Flaggskepparen s:1 har tillagts
fastigheten Bilen 4.
De7 om 9 m' av fastigheten Bilen 4 har ti7lagts samrä7igheten
Flaggskepparen s:1.
Areal: 332 894 m'

- : Se aktang Bilen 4 i 7agf skäp
Del om 96 478 m' av rastigheten Bilen 4 utgör rastigheten
Bile.n IZ.
Areal: 236 416 m2
v G v!

|

Output file 16: Ocrad - Image b

|
rOrtS , , - - - - ij
- te - , - - - - \
- 6,
Fr 7 re 19/8-09: Se akt ang 8i7en 4 i Tagr skäp
/Oel om g36 mz
av rast7gheten 8ilen 4 utgör nybiFdning av
samrä77igheten Stape7badden s:4 fse överenkormP_se i avyttrings -
akt ang de7 om ca 3 ozo m' .
- : J39 546 mZ av rast7gheten 877en 4)

--/..: - ' e ak' ang Bi_en 4 i 7ag , ,käo
Oel om 737 mZ aY rastigheten Bilen 4 utgör Galeasen .
Oe7 om Z zzg mZ .
- : 3J6 5g0 m_ ' a" 'ghe'e" 8i'e" 4 'tg_r _a7easen z_
/

Y G v!

F_rts frgn fasti héten

Fr ä B, s_d7
De7 om 63 mZ . i e" 4 __ä __ - __p - _____

av samfä)igheten F7aggskePparen S:I har ti7lagts
rastigheten 8ilen 4.
/Oel om g mZ
av fast_gheten 8ilen 4 har ti7_agts samrä_ligheten
F7aggskepparen s:1.
- : 33z 894 mZ

Avst re z3/11 - _olo: Se akt ang Bi7en 4 i 7agr skap
|

Oel om 96 478 m'
av rast7gheten B7len 4 utgór rast7gheten
Bile.n IZ.

_: z36 416 mZ

v G v!

Output file 17: *Ocrad – Image c*

: Se aktang 8ilen 4 i 1agf skåp
1 Del om l 017 m' av Bilen 4 utgör Klippern 1.
| Del om 813 m' av Bi7en 4 utgör Klippern z.
l Del om 631 m' av Bilen 4 utgör K7ippern 3.
l Oe .om l 078 m' av Bilen 4 utgör Klippern 4.
f Del om z z55 m' av 8ilen 4 utgör _osterbåten l.
l Del om l 627 m'. . . .
t Del om l 7g4 m2 __ B!'e" 4 "g , Wsterbaten z
Bllen 4 utgör Koggen 1.
| Del om l 666 m' av Bilen 4 utgör Koggen z.
Areal: 2z5 545 m2

_: Se akt ang delar om ca zo60 m2 resp 11 m2 av
fastighéten Bilen 7 i lagf skåp
D-om 759 mZ av Bilen 4 utgör Bilen 13. A 515 och F 11
. Areal: zz4 786 mZ z o o 9

Output file 18: *Ocrad – Image d*

_agtigbeten Oxie 1:5
_r l re 1 4- 9: se akt aAg Pastigh Oxte 1:1 m fl
Det om 10 801,7 a 9amPälligb _ten Oaie s:2 har tilLagtg
Pagttgbeten ie ' : ' g q / ' íg g'

_ rt __ 5p, __ __ __ ..
__ , .

Enl avtal _ har Zel av fastieheten Oxie 1:5 avyttratsl
19 9
SE ak_ an_. c'd T8.ooo l_ av fa_ tigheten Oxie 21;2. Delen har
vid rrgl. reg. 15/1-81 visat sig innehålla T.2g7j9 _.
: 2 227 292,6
En ._, 'trd_ -il 2 / 9_1 Ehr ._l on, __.bl_, y __ TV OyLi_ 1:5,
avs.
t
T._.i I__ fas__S_c-ten E_rnl äf,lni_berl 1_nvT/L.i_a_s. .'_ __ "a'
.Av9t reg l-/3-88: se akt ang Oxie 1:1 m fl
I_el om l_ 225 m_ ay Oxie 1:5 uteör rornl_mningel_ 'L
:_ 2 18o _s2 m

Output file 19: *Ocrad – Image e*

. . bn __Sy
1 .
Fz, stighete_ Oxie 1:5
Bilää_ eAom_ anlägggnN6 _v f_stigh_tErna Oy-ie 1:1,
GxiP, IT:2, Oxi_ 9:id OyL_e 20:8, Oxie 25:9, Oxie _8:v och
O-rtie 29:1.
_ 2.231.22 __
.Särskilt .na_n_ enl.. reg. besl?t. den .1S/_?92: Fr.e.d.?_.ksg_e%hg
Jqr _ re _17_fÖ- r : se akt. rvg . stg 2440 m. fl. i rosie
Dpl_r cIn 85,G ___.v s_nf_l_d_. _mr_E_.n__, 1_956_0 _ _ . __ -
f_bh_ omr___.i _hr___.s:b. ___.o?,_ T% __v_.t%Zsäg_n. 2E__? __ 1
,
_. ___.4 ___.v s___.äL ___.ST,2_6g _z,r. fiV__pcts r__t_C_%t_n
I ___.:

Output file 20: *Ocrad – Image f*

ror_ f_ än
Fzs_ig_eten __ e __ 5 4.
__| om_ Z ___.? ast__ eten_ Ox ___.l_5 utgör __ amräll __ g_ eten_ __ __ __ __
___.t á_7_ A ___. / _Óo6_ ___. - Ón_ z ___. m' a-v . aštig_- té __ __ __ __
. ___.e ___.5 avyttrats __ __ __ __ . __ __ __ __ i ___. - __ __ ! __ __ __ __ .

Output file 21: *Ocrad – Image g*

ror_ f_ än
Fzs_ig_eten __ e __ 5 4.
__| om_ Z ___.? ast__ eten_ Ox ___.l_5 utgör __ amräll __ g_ eten_ __ __ __ __
___.t á_7_ A ___. / _Óo6_ ___. - Ón_ z ___. m' a-v . aštig_- té __ __ __ __
. ___.e ___.5 avyttrats __ __ __ __ . __ __ __ __ i ___. - __ __ ! __ __ __ __ .

Output file 22: *Ocrad – Image h*

t ghPt&z k 4 , , , , , ,
e&n, R m á , , , , , ,
E G, ma a
OW g g, k ma RQr & 2q , ' & p š hRg R,
OW g S &v er 2 , , , , , , , , , , , , , , , ,

Output file 23: *Ocrad – Image i*

Output file 24: *Ocrad – Image j*

Im_ge_1,_hown_n_Appen__A)
 - T_r_Lt 16.97% y5.59% 16.72% y6.7% 16.4%
 - OLR_d 2_.96% 57.9y% y6.2% 75.7% 51.51%
 - Op_nOCR 2y.y4% 49.01% 17.yy% 7.442% 17.6%
 - OCR 62.62% 41.04% 77.6% 62._6% 7_.5y%
 - OCRpu_y2. _7% 15.79% y7.64% 41.52% y9.y%
 TOCR 4.74% 4.41y% ly.y22% 7.207% 11.26%
 AhhW 9.51y% 29.16% 24.56% 27.0y% 15.y1%

T_hl_l: Re,_TILT,_of_the_OCR,_n,.

u

w
u

h
A 4u
h

2u

u
Im_ge 1, hown _n Appen__ A)

& T_r_Lt 11 OLr d 111 Op_nOCR _| _OCR
Oo OCRupu_ || TOCR 11 AhhW

Figur 2: ILLT1, _tr_t_on of the ,___n re ,_TILT,_.
Th_m__n _rrur v_lu_ in inLr____ing urd_r fur __Lh uf th_ OCR tuul_ i_ _huwn
h_luw:

TOCR _169%
Abbyy 21.11%
OpenOCR 22.9 %
Tesseract 24.49%
OCRopus yy.44%
Ocrad 50.07%
GOCR 64.55%

6

Output file 25: *Ocrad – Image k*

Testing OCR tools

Angelica Gabasio

1 English

This is a page with a picture to test the OCr tools.
The page contains text in English and Swedish.

Figure 1: Thi_ is a cat.

2 Swedish

Och __ lite svenska tecken. B_tcn _kcr i __jön.
Det __ är hlött.

1

Output file 26: *Ocrad – Image l*

(3 _iLhL Lh_r_ Ler rzLo___ 'L' _Esn
R__G_L__, h _HPáp_D
z_)_

Output file 27: *Ocrad – Image m*

The algorithm u_c_a database_c_o_. letter_ and numeral_ to match the character_. The algorithm will calculate a matching value_ or each character and collect the option with the bc_t value.111
A way o_. improving the output_i_ to check i_. the _canned word exi_t_ in a dictionary. I_. it doc_n't, it i_ likely that _omc o_. the analyzed letter_ arc wrong, and _omc o_. the letter_ may be changed _o the word matche_ a word in the dictionary. Thi_ i_ done by adding a penalty value to the matching value i_. the analyzed character rc_ult in a word not in the dictionary. Thi_ make it po_iblc _ or another letter to get a better matching value i_. the word with the new letter i_ part o_. the dictionary. The letter_ arc not always changed, _omtimes_ a non-dictionary word get_ a better matching _corc, even with the penalty value, than the one with the changed letter_ that i_ in the dictionary.
Even i_. using a dictionary may help the output accuracy if i_ no guarantee that the correct word will be _ound.111 An example o_. the u_c_o_. a dictionary can be seen in Figure 1.
Other way_ to get a better output i_ to u_c machine learning, where the OCR tool can be trained to recognize more language_ or _ont_, or to u_c the knowledge o_. which character combination_ are more likely to appear. The better the image i_, the more accurate the result get_. I_. the OCR algorithm i_ applied on a horizontal text with very little or no noise and the

lcttcr_ arc wcll _cparated _rom each othr it i_ po_ible to gct a vry
good
rc_ult. But i_. the image i_ blurry , warped or noi_y, or i_. the character_
arc
merged, it i_ likcly that the output gct_ wor_c than it would otherwi_c.ll ,
21

The output accuracy will al_o depnd on the image rc_olution.
OCR _oWwarc i_ generally bettr at rccognizing machinc writing than
handwriting. In machinc writing the _amc character alway_ look_ the _amc,
given that the _amc _..ont i_ u_cd. With handwriting the character_ diKcr a
lot
depnding on who i_ writing , and cvcn i_. the _amc pcr_on ha_ written a tct
the _amc character will bc a little diKrcnt each time it i_ written.
Machine
writing oWcn con_i_t o_. more di_tinct character_ , and _incc handwriting i_
more varicd than machinc writing , it i_ harder _or the algorithm to match
the character_ to it _ databa_c.

3

Output file 28: *Ocrad – Image n*

B.3 CuneiForm

Forts
Fastigheten , Bilen 4 ,
Fr 1 re 19/8-09: Se akt ang Bilen 4 i lagf skäp Del om 836 m av fastigheten
Bilen 4 utgör nybildning av
2
samfälligheten Stapelbädden S:4 (se överenskommelse i avyttringsakt ang del
om ca 3 020 m av fastigheten Bilen 4)
2
Areal : 339 546 m
2
iå
Del om 737 m² av fastigheten Bilen 4 utgor Galeasen 1.
Del om 2 229 m av fastigheten Bilen 4 utgör Galeasen 2
2
Areal: 336 580 m²
Fo~ts från fastigheten
Bilen 4
Fr 1 re 30/7-10: Se akt ang Bilen 4 i lagf skäp
Del om 63 m av samfälligheten Flaggskepparen S: 1 har tillagts
2
fastigheten Bilen 4.
Del om 9 m av fastigheten Bilen 4 har tillagts samfälligheten
2
Flaggskepparen s: 1.
Areal: 332 894 m
2
Avst re 23/11-2010: Se akt ang Bilen 4 i lagf skäp
av fastigheten Bilen 4 utgör fastigheten
2
Bilen 12.
Areal: 236 416 m²

Output file 29: *CuneiForm – Image a*

For ts
Fastigheten Bilen 4 , ,
1
Del om 836 m av fastigheten Bilen 4 utgör nybildning av
2
samfälligheten Stapelbädden S:4 (se överenskommelse i avyttringsakt ang del
om ca 3 020 m av fastigheten Bilen 4)
2
Areal: 339 546 m
2
I
Del om 737 m² av fastigheten Bilen 4 utgor Galeasen 1.
Del om 2 229 m av fastigheten Bilen 4 utgör Galeasen 2
2
Areal: 336 580 m²
Forts från fastigheten
Bilen 4
1
Del om 63 m av samfälligheten Flaggskepparen S: 1 har tillagts
2
fastigheten Bilen 4.
Del om 9 m av fastigheten Bilen 4 har tillagts samfälligheten
2
Flaggskepparen s: 1.
Areal: 332 894 m
2
Avst re 23/11-2010: Se akt ang Bilen 4 i lagf skäp

Del om 96 478 m av fastigheten Bilen 4 utgör fastigheten
2
Bilen 12.
Areal : 236 416 m²
VGV!

Output file 30: *CuneiForm – Image b*

Forts Fastigheten Bile ~ 4 Fr 1 re 19/8-09: Se akt ang Bilen 4 i lagf skåp
Del om 836 m av fastigheten Bilen 4 utgör nybildning av samfälligheten
Stapelbaden S:4 (se överenskommelse i avyttringsakt ang del om ca 3 020
m av fastigheten Bilen 4)
2 Areal: 339 546 m
2
/- Del om 2 229 m av fastigheten Bilen 4 utgör Galeasen 2
2 Areal: 336 580 m²
VGV! För ts frå n fastigheten Bilen 4 Fr 1 re 30/7-10: Se akt ang Bilen 4 i
lagf skåp Del om 63 m av samfälligheten Flaggskepparen S:1 har tillagts
2 fastigheten Bilen 4. Del om 9 m av fastigheten Bilen 4 har tillagts samfä
lligheten
2 Flaggskepparen s:1. Areal: 332 894 m
2 Avst re 23/11-2010: Se akt ang Bilen 4 i lagf skåp Del om 96 478 m av
fastigheten Bilen 4 utgör fastigheten
2 Bi 1 en 12. Areal : 236 416 m²
VGV!

Output file 31: *CuneiForm – Image c*

Enli t avtal A 515 och F 11 har, del om cp 760 m av fastigheten
2
2009
akt ang delar om ca 2 060 m resp 11 m
2
lagf skåp. Avst reg 16/3-2011
Bilen 4 avyttrats. Se
av fastigheten Bilen 7
Enligt avtal A 34/2010 har del om ca 2 400 m av fastigheten
2
Bilen 4 avyttrats. Avst reg .15/7-10
Avst reg 15/7-10: Se avyttringsakt ang del om ca 2 400 m av
Bilen 4.
(Del nm 2 a2B m² av Bilen a utgör fastigheten Klyvaren 1. a 3a/2Dl
Del om 1 3|2 m² av Bilen 4 utgör fastigheten Klyvaren 2.
Areal: 332 840 m²
Forts sid...7
Se akt ang Bilen 4 i lagf skåp
4 utgör Koggen 2.
Areal:
Avst reg 16/3-2011: Se akt ang delar om ca 2060 m² resp 11 m² av
fastigheten Bilen 7 i lagf skåp
Del om 759 m² av Bilen 4 utgör Bilen 13. A 515 och F 1]
Areal: 224 786 m²
200 9
I Del om j Del om j Del om ~ Del.om I Del om ~ D.1 om t Del om 1 Del om
1 017 m av Bilen

2
813 m av Bilen 4
2
631 m av Bilen 4
1 078 m av Bilen
2 255 m av Bilen
2
1 627 m av Bilen
1 784 m av Bilen
1 666 m av Bilen
2
225 545 m2
4 utgör Klippern 1.
utgör Klippern 2.
utgör Klippern 3.
4 utgör Klippern 4.
4 utgor Kosterbåten 1
4 utgör Osterbåten 2.
4 utgör Koggen 1.

Output file 32: *Cuneiform – Image d*

Fastigheten Oxie 1:5
Pr 1 re 1 4-79: se akt ang fastigh Oxie 1: 1 m fl
Del om 10 801,7 8 av samfälligheten Oxie s: 2 har tillagts
fastigheten Oxie 'I: 5 (W 29/197))
~'. X 39 x8,~

A61 260
Enl. avtal j979 har åel av fastigheten Oxie 1:5 avyttrats ,
se akt ang. ca 78.000 uf' av fastigheten Oxie 21;2. Delen har
vid frgl. reg. 15/1-81 visat sig innehålla 7.287,9 m .
areal: 2 227 292.6 9)
Enl. avtal A 20/1981 har del om ~ 2,615,9 rn' av Oxie 1:5 , avs.
att bildas f as+9
t-gheten Pornl ämningen 1 avvttras . (~<å', . -4y')
r eg 17/3-88 : se akt ang Oxie] : 1 m
Bel om 14 225 m ay Oxie 1:5 utgör Fornlämningen 2
2 180 452 m

Output file 33: *Cuneiform – Image e*

Fastigheten Gxie 1: 5
Bildad genom s~anläggning sv fastigheterna Qxie 1:1 ,
Gxie 17: 2 , Oxie 19: 1 , Gxie 20: 8 , Oxie 25: 9 , Oxie 28: 1 och
Oxie 29:1 .
'Areal efter sa~anlä- rin : 2.251.222 rS
Särskilt namn enl. reg. beslut den '15/4-92: Fzedriksbe~g
Fr 1. re ~ 17 10- : se akt ang, stg 2440 m fl. i Posie Delar om 85,0 rE av
samfällda område sdj , 1.956,0 m' av samf ällds området Qxie s : 1 .
"~.507, 2 iP , av stadsägan , 28 4 och 22.122 4. ;E av s+adsägar 2467 ,
2468 hsr tillagts f astigheten
Qx" 1: 5.

Output file 34: *Cuneiform – Image f*

For 's från
Fastigheten Oxie 1:5
Del om 1 261 m av fastigheten Oxie 1:5 utgör samfälligheten

Gånggriften S:1.
Areal: 2 041 038 m

Enligt avtal A 329/2006 har del om ca 23 300 m av fastigheten
2
Oxie 1:5 avyttrats.

Output file 35: *Cuneiform – Image g*

For 's från
Fastigheten Oxie 1:5
Del om 1 261 m av fastigheten Oxie 1:5 utgör samfälligheten
Gånggriften S:1.
Areal: 2 041 038 m

Enligt avtal A 329/2006 har del om ca 23 300 m av fastigheten
2
Oxie 1:5 avyttrats.

Output file 36: *Cuneiform – Image h*

. E@t'\$94 4~\$ A 329/2~ 540 4@3 ~ C4 23 ~ S 4V f@S4\$9W'W5
QX)e 1:5 eyyttree45 ,.

Output file 37: *Cuneiform – Image i*

Forts från
Fastigheten Oxie 1:5
Del om 1 261 m av fastigheten Oxie 1:5 utgör samfälligheten
Gånggriften S:1.
Areal: 2 041 038 m

Enligt avtal A 329/2006 har del om ca 23 300 m av fastigheten
2
Oxie 1:5 avyttrats.

Output file 38: *Cuneiform – Image j*

l,)1)li' 1; li'r i(1 t i/ t/ir (/Cii' rn
llti nti, nt it(in r ilni ut nti ti,i nt:, inili t lin i, lt i)l tlti
((('1(ti)i)l t l nt
lii lint:
'J'OCB ebb> r OpeuOCB 'J'enner tct OCBopian Ocr tcl COCB

1 t , , nti ' , : I 1 1 tint! i(n ii/ t /ir rn rr i(1 t .
. 1 () ' , ') 1.11 , ') 1. 1 ' , ,) . 1 1 ')(1.(1 f) 1. 1 1

Output file 39: *Cuneiform – Image k*

Ti uti o ()C.'H tOOls
Angelica Gabasio 1 English Tlus (s,(p " (u(l h,(ln(I ur(' I() I(st I
h(' OOB louis. Th(p,((uul,(n(s I(sl u(Eu hsh,uul Su((hsh.
Figure D This is a cat. 2 Swedish Och så lite svenska tecken. Båten åker i
sjön.
Det är blött.

Output file 40: *Cuneiform – Image l*

p()<g(4h.-pro,eke~ c >anile -n
VJ

Ail
GE. Ui ~A, C>HEW& i 0

Output file 41: *Cuneiform – Image m*

) ili ', Ii "r) llti1111 IL i', ,I rf, lt ,lf), 1, i' r)i ii 'tti 'I, Illrf 1111111i ,
I IL tr) 111 Itr il tili ' I il 11à ,ii tr I, . ')' llr ,ii , " r)1
itliiii il ill i,ili iil ,itr ,i Il ,iti 11ii , " v, iiiir ir)1 r,ii ll i
il ,il ,ii tr I ,iiiif , I'lr 'i t, till ' r)f)tlr)11 Ll itll till I)r, t v,
tltti .[i]
.) Ll IL r)i 1111I)li)L 111,'~ till ' r)litf)lit I, tr) I ilr 'I il li till ',
I 1111111'rf Lli)lrf 1'Xi, t, 111 I rflr t l))11 11 LI ii It rfr)ii 11 t
. It 1 il)Il 'il til It , r)lli ' r)i tili ' 111 til /I'if ii 't ti I, Ili '
Ll Ii)il ~.,tilif, r)1111 ' r)i till ' ir 't tr 'I, 111,IL I)r' I il,111" I'
rf, r) till ' Lli)lrf 111,1tr ilri ,I Lli)lrf 111 till rlir t i))ii ,ti)I
'1 lli , i, rlr)111 I)y ,trrlrii , " ,I f)r Ii ,tlt y v, tltti tr) till Iil
,ttr llii , " v, tltti Ii iili ' 111 til /I'rf I il 11 ir ii 'I I il Iiii
111 I Lli)I rf 11))i 111 iili ' rflr il))11 11 \.) 111 111 liliit f)r)
,, il)lr ir)1 ,tiir) till I ir ttr I tr),"I t ,i I)r t tr I 111 itr llii
," v iliir ii till Ilr)1 rl Iilitil tili ' Ili Ll ii t ti I 1 I) 11 t r)i
tili ' rflr t l))11 IIL.) ili ' ii t ti I, Ili ' 11))t tiil IL, I il 111~1
if. , r)1111'tlllr ' , I 11))11-rflr t l))11,11 I Lli)I rf "I't , ,I I)r
't tr 'I 111,1tr i1111" , I r)lr ' ltL1'11 Llltil till ' I)i'11 lit l
lili ' . fil 111 tili ' r)lli ' Llltil tili ' I il tll ~i 'if ii 't ti 'I, fil
It 1 111 tili ' rflr tlr)11 II I . L L1'II li il , III," I rlir tir)11 ,il
\ Ill i) Ili 'if tili ' r)lit f)lit ii i ill ii \ it i, llr) " il ,il
illr 'I , tll ,it till r r)1 11 r t Ilr)1 il 11 ill I)r ir)tilrl.[1] . 41
I x, tiifl)li r)i till ii , I r)i ,i rlir t i))ii ,ti y
ill I)r', I'I'II ill i'i illr ' 1.
() till I Ll IL , tr) "I t I I)r t tr I r)lit f)lit I, tr) il I' Ill ir
lllllr ' lr ' tl lllll ~. Krill lr ' till ' (((l) tr)))l i, ill I)r ' tl
illlr ' il tr) lr 'i)), " 111/I' I11r)lr ' i ill , " il I, " ri r)I ir)lit , .))I
tr) il , I' till ' lill))till ' il , " I' r)i Llllii ll I il il ir tr I I r)
lill)lll itl))ll , il r' I11r)lr ' lill ' i) tr) if)f)r ' Il.
) ili ' I)i 't ti 'I tili ' 1111,1" i ' I, . tili ' 111))li ' ,Ir I Iil ,lti ' tili '
Ii ' , Iiit " i 't , . ii tili ' (((i) tl , ' ~i)11t1111 I, if)f)llr 'rl r)11
I llr)1 I/r)lit Il tl 'x t Ll ltll L1'I I lit tlr ' r)I llr)1 , I'
illrl till ' lr t tr I ilr ' Kil il I f) il itr rl il))111 I' ir ll r)till
I lt I f)r), , 11)lr ' ti) , ' ~1 t I L1 I I , ' ~i)i)il lr ' , illt . f fttt li
till ' lill I" I' I, I)lttl I I . Ll il f)r 'rl r)I llr)1 , LI r)I li till

' r ll il ir tr 'I, il r' 1111 I "I il. it i, lilll ly t il it till r)
 iitf) iit "I t, ilr) 1 ~1 till tii it ilr) ttilil r) till I il i, I .~1. 2~ 1
 lll ' r) lit f) lit , ii i Ill, ii \ Krill , tl , r) ilr 'f) r' llil r) ll till ,
 lill I "I, lri r) illt1)) ll.
 ()(1) r) ittl il r' I , "I llr I Ill) I)r'ttr I it lr r r), " lll/ lll , " Ill
 Ir lllllr , Lll itiii , " t il tii ll illrftl lltlll ~. lll Ill ir lllllr ,
 LII ltlll ~ till , illlr , r ll il ir tr 'I tltl IL, lr)r) ll , till , illlr
 ' . , ' 111'11 til It tili , 1111i' ir) lit 1 Ii , i 'rf. A Itii il lllrftli
 ltill ~ tili ' I il 11 Ir ti I, rfliii I I ir)t rlr f) l llllll ~ r) ll
 Lill)) I, Ll lltlll ~. Illrl ltLlll li till , illlr , f)r I, r) ll ll I,
 Llllt tr II I tl xt. tili ' , 1111i' I il , ll , lr ti I Lllii I)i , I ilt tii ,
 rfliii Ii ' Ilt i , Ir il tllli ' It 1 LII ltti ' 11. L i , tr 111111 ' LII
 ltlil ~ r) iii ' 11 i i) 11, I, i r)i 111)) li ' ifl , illli i i il 11 ii ii ' I,
 . Illif , Illi i ' il tili ftl lltliil " I, Ilr) I I II ti il il t il , tii
 111, I 1 lliiil il I i till , " . i t i , ll , ti 1 ll I ir) I till , tl , "
 r) I i tillit t r) 111, I t 1 ll till ' r ll il ir tr 'I, tr) lt , rl it
 il), i , I .

Output file 42: *Cuneiform – Image n*

B.4 GOCR

--D DD_eeeol_r_tooom_mm296_3r2_m_a2_m9m - 99 - 99 t
 rGa_e_asen/_2____a9ts
 fastigheter1 _ilen 4 -

Fr 1 19/8 09: Se akt ang Bilen 4 i 1agf skåp
 2 av fast_.gheten Bl.1en 4 utgo..r nyb_.1dn_.
 samfälligheten Stape1bädden S:4 (se överenskō nme lse i avyttrings-
 2 av fast_. heten Bl.
 2

A_t r 30/9 09: Se akt ang Bilen 4 i 1agf skåp
 Del om J3J m2 av fastigheten Bjlen Q utgör Galeasen 1
 2 av fast_. eten B_.len 4 utgo..
 Areal; 336 580 m2

V G V!

8, fastl_ _ete_ Sid J
 Bilen 4
 fr 1 r 30/7 10: Se akt ang Bilen 4 i 1agf skåp
 2 av samf.a.111. heten F_a skepparen s.1 har t_,
 fastigheten Bilen 4.
 2a
 V faStl9hete_ 811en 4 har tl11a9tS Samfäll119heten
 flaggskepparen s:1.
 2

A st r 23/11 2010: Se akt ang Bilen Q i 1agf skåp
 De_ om g6 4Jg m2 av fast_.ghete, B_tle, 4 ut o.., f,st_t
 BjTen 12.
 Area; 236 Q16 m2
 V G V!

Output file 43: *GOCR – Image a*

--)____ t____D_ADDerele_ao_ommt_ 236332
 2 2____m989m4avmasvamffta as_t____99hheetee_nnBF____lae9n9s
 Qkeuptp_9ao_r FenGas____elashean_r_r2t____a9ts ,
 f_r_t ,,-, - --
 f_tightet_ _i_en 4

fr 1 19/8 09: Se akt ang Bilen 4 i 1agf skåp
 2 samfälligheten Stape1bädden S:Q (se överensko nme lse i avyttrings-
 2
 2

A vst re 30/9 09: Se akt ang Bilen 4 i 1agf skåp
 De_ om J3J m2 ay fastjgheten Bilen 4 utgöF Galeasen 1
 2

Areal: 336 580 m2

/

v G V!

,
 Bielen 4 -rts _rán _ast_g-hetgn S_d J
 fr1 30/7 10: Se aktang Bielen 4i 1agfskåp
 2
 fastigheten 8ilen 4,
 2
 el Om 9 m aV faStl9heten B1le, 4 haF tl11a9tS Samfall19heten
 flaggskepparen s:1.
 2

Avst re 23/11 2010; Se akt ang Bielen Q i 1agf skåp
 0el om g6 4Jg m2 av fast_g_eten 8_le, _ ut o.., f,,t__
 8ile.n 12
 Areal; 236 416 m2
 V G V.!

Output file 44: GOCR – Image b

_____ 9 t9o__9F9Ga7peapsenl_r9_a9ts
 fas_rt_i g h e t e J1 _ i _ e _ _ 4
 Fr01 re_ 1 9 / 8- O 9: S e a k t a n g _ i 7 e n 4 i 1 a g f s k å p
 0e7 ,m 8 3 6 m2 a v f a s t _ g h e t e n g i _ e n 4 _ t g ö r n y b_ J d n
 n a v
 samfä7li g h e t e n S t a p e 1 b _ d d e n S: 9 (s e ö v e r e n s _ o
 nme l s e i a v y t t r i n g s
 a_t a n g _ d e _ o m c a 3 o 2 o m 2 a v f a s t i g h e t e n g i 7 e n q)
 A,e a 1 .. 3 3 g 5 4 6 m 2

Avst re_ 3 O / 9- O 9: S e a k t a n g 8 i 1 e n 4 j 1 a g f s k å;
 De_ om 7 3 7-, m 2 a y f a s t . g h e t e n B _ e n q u
 De7 om 2 2 2 g m2 a v f a s t i g h e t e n B i _ e n g u t ö_, G a 7 e a s
 e n 2
 Area7: 336 5 8 O m 2 / V G V !

Fort _rán f a s t _ g h e t e n S i d J
 8ilen 4
 fr07 re0 30 / 7- 1 O: S e a k t a n g 8 j 7 e n 4 j 1 a g f s _ å
 Del om 63 m2 a v s a m f a _ _ 7 7 g h e t , n F 7 a g g s k e p p a, e n s..
 1 h a, t _ 7 7
 fastigheten 8 i _ 7 e n 4.
 De_ om g m2 a, f a s t i h e t e n B 7 _ e n q h a, t 7 l _ a t s, a m f a _
 7 7 . h e t e n
 Flaggs_ep par e n s: 1 ,
 A,ea7.. 3 3 2 g g 4 m2

Avst re0 23/11 - 2 O 1 O: S e a _ t a n g B j 1 e n 4 j 1 a g f s k a '
 De7 om g6 47g m2 a, f a s t _ _ g h e t e n g _ _ e n q _
 8ile.n 12.
 Area7: 236 4 1 6 m2 V G V. !

Output file 45: GOCR – Image c

n t 2 / 1 _ : Se akt ang Bi_en 4 i 1 agr _ ká_
2 av B_en Q utgo.. F K1_.
2 av B_en - utgo.. r K_.
2 av 8_1en _ utgo.. r K1_.
2 av B_len Q utg.o.r Kl_.
2 av B_len 4 utgo.. r _osterbao
2
Del om 1 627 m av Bilen utgör sterbåten 2,
Del om 1 78Q m av 8ilen 4 utgör Koggen 1.
De_ om 1 666 m2 av 8___en 4 utgo_., Koggen 2
Are,_.. 225 545 m2
Avst reg 16/3-2011: Se akt ang delar om ca 2060 m2 resp 11 m2 av
fastigheten Bilen 7 i 1agf skåp
D'1 om 759 m2 av Bilen _ utgör Bilen 13. A 515 och f 11
Areal: 224 786 m2

Output file 46: *GOCR – Image* d

n o 2 Tt r22 7 g 1 62l29o /t 665o e2 1 t o t e
 o l rtJ 61 r c Jt l r l 3v o rv t avst)
 &stig beten O ie ; 5
 0 e 0 1 0 4 7g; ge abt 8Dg 8stigh Ox te 1 : 1 - - -
 ae 1 beten o ' e t: 5 (A 22 /1 , - ____)
 !f i , i 7 ; ,
 E aY g jg har el aV fa5 i he e OX e 1 5 aYy tr t5,
 s k a . ca 78 OO av fa5 tigheten Ox i 2 lj 2 De e ha
 vid f gl eg. 15/1-81 visat sig in eh å la 7. 2 7 9 -

Output file 47: *GOCR – Image e*

_____p_y_e_ar_m_g_vör_s_t_n_2_&
 _1_3_nrTT_f_t_l_lyln_ID_tr
 _2_2_ttol_2_t_2_2e_t_1_9_o_r_17_6_?_m0_y
 'astigh_e_e_u_xje
 gi_a_a_e_,_2n_ägg_i_wv f3s_t ig he te n2_0_je
 xie?_2t_Xi_9,_it_i_e 2_8t_0X_e_25_9t_OXiE_8_C
 _j_29:_.
 Särskilt na_mm_enl. reg. beslut_d_en 1_5/4-92: f_ed_r_sb._e_r.
 ,_o_e_0 io-' : se akt, _g_s_g_4_?_i_F_D
 f_j_om_å_?_. = ., Z,_4,_7?/_r_. ; t_säe_n_2_?
 J_i_?x_4,-_?_, c_v_t,n_,'_ J_6?._,_6_5_J
 ti_n_t_3,5_icn_,t_?5_i_5,
 ?5_i_5,

Output file 48: *GOCR – Image f*

— AE _rea 2 o 4 l o 38 m — — — — —
or f án
f2stiy_ , n _e i:5 4.
el om l 261 m av ^ astighet_en_Oxie l: 5 u_t_gör sam_fäll gheten
Gång_griften S :l.
2 /

n119t aVtal A 329 2006 haf del Om_Ca 23 300 m aV faStl9_heten
oxie l: 5 avytt ats.

Output file 49: *GOCR – Image g*

— AE _rea_ 2 o_ 4 l_ o_ 38_ m_ — — — — —
or_ f_ án
f2stiy_ , n_ e_ i:5 4.
el_ om_ l_ 261_ m_ av_ ^_astighet_ en_ Oxie_ 1:_ 5_ u_ t_ gör_ sam_ fäl1_ gheten
Gång_ griften_ S_ :1.
2 /

Output file 50: *GOCR – Image h*

Output file 51: *GOCR – Image i*

— AE — m e a 9 2 o 4 1 038 — m —
— a v f 00 t 9 he ten t
or f án
f2stiy e+ en _ e i:5 Q.
e1 om l 261 m av ^ astighet_en_Dxi_e 1:_5 u_t_g_ör sam_fäll ghet_en
G ång_gr_i f_ten _S _:1.
2
- - - - , ' ^ ' - - , . - -
. t avta A 32g/2006 har de_ om ca 23 300 m 2 —
o_x e l :_5 _a vy_t t_a ts ._- — — — ,
?
>
—
—

Output file 52: *GOCR – Image j*

$$4 \text{---} qJ \text{---} u \text{---}$$

JTnnge (,sh,own n Appen_zi' AJ
 h c d ,
 T,ss,r_ct 16.97%o 35.59%o 16.72%o 36.7%o 16.4%o
 Ocr_d 2.96%o 57.93%o 36.2%o 75.7%o 51.51%o
 op, nocR 23.43%o 49.01%o 17.33%o 7.442%o 17.6%o
 m
 OCRupus' 32.7%o 15.79%o 37.64%o 41.52%o 39.3%o
 TOCR 4.74%o 4.413%o 13.22%o 7.207%o 11.26%o
 Ahhvyy 9.513%o 29.16%o 24.56%o 27.03%o 15.31%o

T_hl, l: ne, s1_lt, s_oi_th, e_OCn, scnn, s.

YU

$$-\mathbf{q} \quad -\mathbf{U}$$

—

4u

—

2U

M h C d ,
JTnnge (, sh, own n Appen_zi' AJ

T, s 's ', r _ ct l _ Ocr _ d Il _ Op,nOCR GOCR
— — OCRupus' I _ TOCR Ahhyy

Figure 2: Jll1_struct on oil theft, screen re, s1_struct, s.

Th, m,_n ,rrur v_lu, in incr,_s'ing urd,r fur ,_ch uf th, OCR tuuls ' is ' s' huwn h,luw:

TocR	16.9%
Abbyy	21.11%
OpenOCR	22.9%
Tesseract	24.49%
OCROpus	33.44%
Ocrad	50.07%
Gocr	64.55%

6

Output file 53: *GOCR – Image k*

-4 ?v?v x? 0 ? 00 y u? t c? 0 t? ? J 0 c ?
-? -t0 -Jh? -lr / r t c0tf r f -) -

Testin_ OCR tools

— —, ; 0 c 0
—, J —
—, ; , , ; .
—, m.
— 0, e, , 0, , ? ' =,
?, , 0

Fig. 1: This is a c_t.

2 Swedish
Och så litc svns_ techcn. Båtcn åhcr i sjön.
Dct _r hlött.

1

Output file 54: *GOCR – Image* l

_r_a_ _ _ _ _ \ _h_ _ _ _ _ \ _u_ _ _ _ _ e_ _ _ _ _ e_ _ _ _ _ u_ _ _ _ _ _ \ _i_ _ _ _ _ n_

Output file 55: *GOCR – Image* m

Tl_c algoritl_____ uscs a datal_asc of' lcttcrs a_d _u__cralrs to ___atcl_ tl_c cl_ar-

actcrs. Tl_c algoritl_____ will calculate a ___atcl_i__g valuc f'or cacl_ cl_aracter a_d

sclects tl_c o__tio__ witl_ tl_c l_cst valuc. 1

A wa_,T of, i____rovi__g tl_c out__ut is to cl_cch if' tl_c sca____cd word exists i__

a dictio__ar_,T. If' it docs__t., it is lihcl_,T tl_at so__c of' tl_c a__al_,T cd lcttcrs arc

wro_g., a_d so__c of' tl_c lcttcrs __a_,T l_c cl_a__gcd so tl_c word atcl_cs a word i__

tl_c dictio__ar_,T. Tl_is is do__c l__,T addi__g a __c__alt_,T valuc to tl_c atcl_i__g valuc

if' tl_c a__al_,T cd cl_aracter rcsult i__ a word __ot i__ tl_c dictio__ar_, T. Tl_is ahcs

it _ossil_lcf'or a__otl_cr lcttcr to gct a l_cttcr ___atcl_i__g valuc if'

tl_c word witl_

tl_c cw lcttcr is__art of' tl_c dictio__ar_,T. Tl_c lcttcrs arc __ot alwa_,Ts cl_a__gcd.,

so__cti__cs a__o__dictio__ar_,T word gcts a l_cttcr ___atcl_i__g scor.,

cvc witl_ tl_c __c__alt_,T valuc., tl_a_ tl_c o__c witl_ tl_c cl_a__gcd lcttcrs tl_at is i__ tl_c dictio__ar_,T.

Evc__if' usi__g a dictio__ar_,T __a_,T l_cl__ tl_c out__ut accurac_,T it is o_guara__tcc

tl_at tl_c correct word will l_c f'ou__d. 1 A__ cxa____lc of' tl_c usc of'

a dictio__ar_,T

ca_l_c scc_i__ Figure 1.

Otl_cr wa_,Ts to gct a l_cttcr out__ut is to usc ___acl_i__c lcar__i__g., wl_crc tl_c

OCR tool ca_l_c trai__cd to rccog__i_c ___orc la__guages or f'o__ts., or to usc tl_c

h__owlcdge of' wl_icl_ cl_aracter co__l_i__atio__s arc ___orc lihcl_,T to a car.

Tl_c l_ctter tl_c i_agc is., tl_c ___orc accurate tl_c rcsult gcts. If
 tl_c OCR
 algoritl_ is a_lid o_a l_ori_o_tal tctx witl_vcr_,T littlc or __o
 _oisc a_d tl_c
 lctters arc wcll sc_arated f'ro__ cacl_ otl_cr it is __ossil_lc to gct a
 vcr_,T good
 rcsult. But if' tl_c i_agc is l_lurr_,T., war_cd or __ois_,T., or if'
 tl_c cl_aracters arc
 ___rged., it is lihcl_,T tl_at tl_c out_ut gcts worsc tl_a__ it would
 otl_crwisc_1., 2
 Tl_c out_ut accurac_,T will also dc_c_d o tl_c i_agc rcsolutio_.
 OCR sof'tware is gc_crall_,T l_ctter at rccog_i_i_g ___acl_i_c writi_g
 tl_a
 l_a_dwriti_g. I___ acl_i_c writi_g tl_c sa_c cl_aracter alwa_,Ts
 loohs tl_c sa_c,
 give tl_at tl_c sa_c f'o_t is uscd. Witl_ l_a_dwriti_g tl_c
 cl_aracters di_r a lot
 dc_c_di_g o_wl_o is writi_g., a_d cvc_if' tl_c sa_c crso_l_as
 writtc_a_tctx.,
 tl_c sa_c cl_aracter will l_c a littlc di_rc_t cacl_t i_c it is
 writtc_. /1acl_i_c
 writi_g of'tc co_sist of' ___orc disti_ct cl_aracters., a_d si_cc
 l_a_dwriti_g is
 ___orc varicd tl_a___ acl_i_c writi_g., it is l_aracr f'or tl_c
 algoritl_to_atcl
 tl_c cl_aracters to its datal_asc.

B.5 OCropus

i.e..
Ports
Fastigheten . Bilen . 4
Fr 1 re 19/8-09: Se azt ang Bilen 4 i lagf skJp
Del om 836 mg av fastigheten Bilen 4 egdr nybildning av samfAlligheten Stapelbddd S:4 (se dverenskoelse i avyttrings akt ang del om ca 3 020 mgav fastigheten Bilen 4)
Areal: 339 546 m²
Avst re 30/9-09: .Se azt ang Bilen 4 i lagf skJp
Del om 73im2 av fastigheten Bilen 4 utgdr Galeasen h
Del om 2 229 m² av fastigheten Bilen 4 utgdr Galeasen k.
Areal: 336 580 m²
Forts fr5n fastigheten
Bilen 4
Fr 1 re 30/710: Se azt ang Bilen 4 i lagf skJp
ys VI
Sil 1
Del om 63 m² av samftilligheten Flaggskepparen S:1 har tillagts fastigheten Bilen 4.
Del om 9 m² av fastigheten Bilen 4 har tillagts samfAlligheten Flaggskepparen s:1.
Areal: 332 894 m²
,Avst re 23/11-2010: Se akt ang Bilen 4 i lagf skJp
Del om 96 478 m² av fastigheten Bilen 4 utgor fastigheten Bilen 12.
Areal: 236 416 m²
Ys vl

Output file 57: *OCropus – Image a*

-'y.,
PortS.
Fastigheten . i1n 4 e
Fr 1 re 19/8-09: Se azt ang Bilen 4 i lagf skJp
Del om 836 m² av fastigheten Bilen 4 utgdr nybildning av samfAlligheten Stapelbddd S:4 (se dverenskoelse i avyttrings akt ang del om ca 3 020 mgav fastigheten Bilen 4)
Areal: 339 546 m²
Awst re 30/909: e Se azt ang Bilen 4 i lagf skJp
Del om 731 sm² av fastigheten Bilen 4 utgdr Galeasen L
lDel om 2 229 m² av fastigheten Bilen 4 utgdr Galeasen 2.
Areal: 336 580 m²
Forts frar.astigheten
Bilen 4
Fr 1 re 30/710: Se azt ang Bilen 4 i lagf skJp
ys Vt
Sif 1
Del om 63 m² av samfAlligheten Flaggskepparen S:1 har tillagts fastigheten Bilen 4.
Del om 9 m² av fastigheten Bilen 4 har tillagts samfelligheten Flaggskepparen s:1.
Areal: 332 894 d2
,Avst re 23/11-2010: Se akt ang Bilen 4 i lagf skap
Del om 96 478 m² av fastigheten Bilen 4 utgor fastigheten Bilen 12.
Areal: 236 416 m²
VG V1

Output file 58: *OCropus – Image b*

```

portS .
renen-,v
Fastigheten ilene 4 s
Fr 1 re 19/809: Se akt ang Bilen 4 i lagomp
a1 ds b5 m2 a, fastigheten Silen 4 agdr nybildnin0 z
% samralligheten StaoNhadden S:4 (se overenskdgmise i aortrings -
okt ang del om ca 3 020 m2 av fastigheten Bilef)
al: 339 546 m2
,Aat re 30/9-09; Se aM ang Bilen 4 i lagimp
61 om 737PA ay fastigheten Bilen 4 utgar sa0easen L
101 om 2 229 m2 av fastijkten 8ilen 4 utsyar saleasen p
st [: 336 S80 m2
Forts fran fastigheten
Bilen 4
y6 yI
Sii 1
Fr 1 re 30/7-10: Se akt ang Bilen 4 i lagoep
el om as md av samalligheten Fla99skepparen S:1 bar tillas
fastigheten Bilen 4.
ai oe s Nz av taytisaetee stien 4 nar t1iaazts samsuiUkus
[Flaggskepparen s:L
aCl: 332 894 m2
AEt re 23/11-2010: Se akt ang Bilen 4 i lagiep
a1 ve 96 Os as raytIonetan stien 4 utner rasttsoten
Bilen 12.
S%(: 236 416 m2
yG VI

```

Output file 59: *OCRopus – Image c*

Enltaat A 515 och F 11 har del om ca 760 m2 av fastigheten
2009
e5ilen navyttrets,Se azt ang delar dom ca 2 . 060 e.resp 11 m2
av d efashighetem Bilen 7d.i slagf skap. Avst reg 16/3-2011
Enligt avtal A 34/2010 e har d del om ca 2 400 mg av fastigheten ee,
Bilen 4 avyttrats...Avst reg 15/7-10 - -
.Avst reg 15/7-1O: Se avyttingsakt dang del om ca 2 400 m2 av
Bilen 4- -
1 Del om 2428 m2 av Bilen s4 sutgbr efastigheteaKlyvaren 1. A34/201
Del om 1 312 m24 av Bilen 4 utgor fastigheten Klyvaren g.
Areal: 332 840 o2
----FOrtS Sid)
Avst re 25/12011: Se akt ang Bilen 4 i lagf skap
Del om 1 017 m2 av Bilen 4 utgor Klipern 1.
Del om 813 mg av Bilen 4 utg6r Klipern d
'Del om 631 mg av Bilen 4 utgdr Klipern l
'Del om 1 08 m2 av Bilen 4 utgdr Klipern 4.
Del om 2 255 mmav Bilen 4 utgdr yosterbaten 1.
Del om 1 627 meav Bi len 4 utgdr Msterbaten 2.
1 Del om 1 784 mc av Bilen 4 utgor Koggen E
[Del om 1 666 m2 av Bilen 4 utgdr Koggen d
Areal: 225 545 m2
Avst reg 16/3-2011: Se azt ang delar om ca 2060 m2 resp 11 m2 av
fastigheten Bilen 7 i lagf skJp
Del om 759 m2 av Bilen 4 utgor Bilen 13. A 515 och F 11
Areal: 224 786 m2
200 9

Output file 60: *OCRopus – Image d*

Faetigheten oxie 1:5

Pr 1 re 1 4-79: ee akt ang fastigh Oxie 1:1 m fl
 Lel om 1O 801o7 3 av eamfalligh 'ten Oxie e:2 har tillagte
 5/107oY
 faetigheten Oxie 1:5 (F st /; rtra /
 4 dt / Abod
 re ,e 8ca.7 -d)
 Enl,' avtaj d bar del av fastigheten oxie 1:5 avybtrans,
 1979
 se akt ang. ca 75.00C vi av fastigheten Cxie 21:2. Delen nar
 vi ,frgl. reg. l51O1 visAt oig inneholla 7.2O7.9 d.
 j: 2 227 292,6 j
 Ems. e7tol A 2Oto5l oaU del om i2.62,oc nan av CZie lt,, ave.
 dh- bolaal fanieHnen por, amningen ,av. crats..c, sh e'C
 Avst Teg 17 /3884 se ak t ang Oxie 1:1 m fl
 Uel oml0 225 me ag Oxiel:5 HtGUR FOtridmoInGen2
 : 2 18O 452 m'

Output file 61: *OCRopus – Image e*

Faiigheten Cxie 1:5
 Li /,--N
 Bildad genom eammanlaggnimg aefaeligheterne Crie 1:1,
 Oxie lac ,,, Oxie 19:1, Oxie 20:8, Oxie 25:9, Oxie 2S4 1 och
 Oxie 29:1.
 'Areaieftersainld-nin:i 2.251.222 3
 Sdrskot ensip deol. reg. beslot den 54/4-92: eedrues049
 Fr 1. re .17 t0-
 i ee akt ang. Atg 244C m.f2. i Foeie
 Pelar om 853C 3 av samfdldn omradet adj, 1.956oC Pn av sam
 fall5e omrdiet Chle s:1, 33.C7o 9 U' Rv etrdshgan . 28e och
 -r.repon;n ev stadengam 7467,246O her tillagte faetighoten
 orie 1:5.

Output file 62: *OCRopus – Image f*

corts fran
 Fastiheten Osie 1:5
 Del om 1 261 m av fastigheten Oxie 1:5 utgUr samfalligheten
 Ganggriften S:1.
 Areal: 2 041 038 m2
 Enligt stdla3k9/206 hanelom ca23 d300 smg ave fastighetenses
 c)
 10xie 1:5 avytrrats.

Output file 63: *OCRopus – Image g*

corts fran
 Fastiheten Osie 1:5
 Del om 1 261 m av fastigheten Oxie 1:5 utgUr samfalligheten
 Ganggriften S:1.
 Areal: 2 041 038 m2
 Enligt stdla3k9/206 hanelom ca23 d300 smg ave fastighetenses
 c)
 10xie 1:5 avytrrats.

Output file 64: *OCRopus – Image h*

, i -- a1
cm
, if No, v.
g -. .
r
. .
,
x
4 o
sed ' oo , O re obe ees. 0 asAse za4od sessoe e wseeisid 0 ses 0 ee , Ma soo - eg y
. o .
W
00
A
. .
'd
L. s
, m
ee
e
t
\$ o a 00
. n
e
E ,
vo +
:
e
'Ze v
,
Y'6
ii
. ,
hvt
e e
r d
s
nAs
sy
s
. .
i1J
o roo
00
e
Soc
i sA
serf ,
P
Ps
en
e
4
e
00
f
s
e
-
4 i . /
= -
-
,
4
e

```
4 1
s
o
r -U
ats
ow
9
e
o
e
A
o
E
.6
isA
```

Output file 65: *OCRopus – Image i*

```
cortsfrah
Fastigheten Osie I:5
Del om 1 261 m av fastigheten Oxie 1:5 utgor samfalligheten
Ganggriften S:1-
Areal: 2 041 038 m2
e e e
i
[EnligtavtalAR3I3006 hanalomca23 300 m2 av .ftighetenesses
j
108ies 1.5 avyttrate , ,
Ins --
eth
AP ---
seg
-i --
a
```

Output file 66: *OCRopus – Image j*

```
Image (shown in Appnds A7
a b c d e
gTact 16.97% 35.59% 16.72% 36.78% 16.49
D Ocrad 28.0% 57.93% 362% 75.78% 51.519
5 OpetCR 23.43% 4901% 17.33% 7.442% 17.68%
, GOCR 62 62% 41 04% 7768% 6286% 7853%
OCRops 32.87% 15.79% 3764% 41.52% 39.38W
TOCR 474% 4413% 1322% 7207% 11269
Abbyy 9.513% 2916% 2466% 27.03% 15.31O
Table 1: Results of the OCR scans
ab
d e
Image (shown in Appnds A7
As Tessnao ls OOw tz OpeuOCR la GOCR
14 OCRopus if TOCR Abbyy
Figure 2: Illustraoon of the scan results
The meau ero value in increasing order for each of the OCR tools is showu
below
TOCR 81699
Abbyy 21.11A
OpenOCR 22989
Tesseract 24 499
OCRopus 33.449
Ocrad 50 070
```

GOCR 64559

Output file 67: *OCRopus – Image k*

```
1 English
Testing OClools
Angehca Gabaslo
This is a page with a picture to test the OCR tools
The page contains text in Enghsh and Swedish
2 Swedish
O s
es
m
f .-
- M- N
'ure l: This is a cat
Och sa hte sveuso teckeu. Bateu aker icou
Det ar blott
```

Output file 68: *OCRopus – Image l*

```
'pfica (Cheir.c cee., orooo
(hN'thetrostAl
,
```

Output file 69: *OCRopus – Image m*

the algorithm uses a database of letters and numerals to match the characters. The algorithm will calculate a matching value for each character and selects the option with the best value[1]. A way of improving the output is to check if the scanned word exists in a dictionary. If it doesn't, it is likely that some of the analyzed letters are wrong, and some of the letters may be changed so the word matches a word in the dictionary. This is done by adding a penalty value to the matching value if the analyzed character result in a word not in the dictionary. This makes it possible for another letter to get a better matching value if the word with the new letter is part of the dictionary. The letters are not always changed sometimes a non-dictionary word gets a better matching score, even with the penality value, than the one with the changed letters that is in the dictionary. Even if using a dictionary may help the output accuracy it is no guarantee that the correct word will be found[1]. An example of the use of a dictionary can be seen in Figure 1.

Other ways to get a better output is to use machine learning, where the OCR tool can be trained to recognize more languages or fonts, or to use the knowledge of which character combinations are more likely to appear. The better the image is, the more accurate the result gets. If the OCR algorithm is applied on a horizontal text with very little or no noise and the letters are well separated from each other it is possible to get a very good result. But if the image is blurry, warped or noisy, or if the characters are merged, it is likely that the output gets worse than it would otherwise[1]. The output accuracy will also depend on the image resolution.

OCR software is generally better at recognizing machine writing than handwriting. In machine writing the same character always looks the same given that the same font is used. With handwriting the characters differ a lot depending on who is writing, and even if the same person has written a lot

the same character will be a little different each time it is written. Machine writing often consists of more distinct characters, and since handwriting is more varied than machine writing, it is harder for the algorithm to match the characters to its database.

Output file 70: *OCRopus – Image n*

B.6 TOCR

Forts 6.

Fastigheten Bilen 4

Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp

Del om 836 m² av fastigheten Bilen 4 utgör nybildning av samfälligheten Stapelbädden S:4 (se överenskommelse i avyttrings-akt ang del om ca 3 020 m² av fastigheten Bilen 4)

Areal : 339 546 m²

Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp

Del om 737 m² av fastigheten Bilen 4 utgör Galeasen 1.

Del om 2 229 m² av fastigheten Bilen 4 utgör Galeasen 2.

Areal: 336 580 m²

V G V!

Forts från fastigheten Sid 7

Bilen 4

Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp

Del om 63 m² av samfälligheten Flaggskepparen S:1 har tillagts fastigheten Bilen 4.

Del om 9 m av fastigheten Bilen 4 har tillagts samfälligheten

Flaggskepparen s:1.

Areal : 332 894 m²

Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp

Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten

Bilen 12.

Areal : 236 416 m²

V G V!

Output file 71: *TOCR – Image a*

Forts 6.

Fastigheten Bilen 4

Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp

Del om 836 m² av fastigheten Bilen 4 utgör nybildning av samfälligheten Stapelbädden S:4 (se överenskommelse i avyttrings-akt ang del om ca 3 020 m² av fastigheten Bilen 4)

Areal : 339 546 m²

Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp

Del om 737 m² av fastigheten Bilen 4 utgör Galeasen 1.

Del om 2 229 m² av fastigheten Bilen 4 utgör Galeasen 2

Areal: 336 580 m²

V G V!

Forts från fastigheten Sid 7

Bilen 4

Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp

Del om 63 m² av samfälligheten Flaggskepparen S:1 har tillagts fastigheten Bilen 4.

Del om 9 m² av fastigheten Bilen 4 har tillagts samfälligheten

Flaggskepparen s:1.

Areal : 332 894 m²

Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp

Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten

Bilen 12.

Areal : 236 416 m²

V G V!

Output file 72: *TOCR – Image b*

Forts 6.

Fastigheten Bilen 4

Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp

Del om 836 m² av fastigheten Bilen 4 utgör nybildning av samfälligheten Stapelbädden S:4 (se överenskommelse i avyttrings-akt ang del om ca 3 020 m² av fastigheten Bilen 4)
 Areal : 339 546 m²
 Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
 Del om 737 m² av fastigheten Bilen 4 utgör Galeasen 1.
 Del om 2 229 m² av fastigheten Bilen 4 utgör Galeasen 2.
 Areal: 336 580 m²
 V G Y!
 Forts från fastigheten
 Bilen 4
 Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
 Del om 63 m² av samfälligheten Flaggskepparen S:1 har tillagts fastigheten Bilen 4.
 Del om 9 m² av fastigheten Bilen 4 har tillagts samfälligheten Flaggskepparen s:1.
 Areal: 332 894 m²
 Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
 Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten Bilen 12.
 Areal : 236 416 m²
 V G V!

Output file 73: TOCR – Image c

Enligt avtal A 515 och F 11 har del om ca 760 m² av fastigheten 2009
 Bilen 4 avyttrats. Se akt ang delar om ca 2 060 m² resp 11 m² av fastigheten Bilen 7 i lagf skåp. Avst reg 16/3-2011
 Enligt avtal A 34/2010 har del om ca 2 400 m² av fastigheten Bilen 4 avyttrats. Avst reg 15/7-10
 Avst reg 15/7-10: Se avyttringsakt ang del om ca 2 400 m av Bilen 4.
 Del om 2 428 m² av Bilen 4 utgör fastigheten Klyvaren 1. A 34/201
 Del om 1 312 m² av Bilen 4 utgör fastigheten Klyvaren 2.
 Areal : 332 840 m²
 Forts sid 7
 Avst reg 25/1-2011: Se akt ang Bilen 4 i lagf skåp
 Del om 1 017 m² av Bilen 4 utgör Klippan 1.
 Del om 813 m² av Bilen 4 utgör Klippan 2.
 Del om 631 m² av Bilen 4 utgör Klippan 3.
 Del om 1 078 m² av Bilen 4 utgör Klippan 4.
 Del om 2 255 m² av Bilen 4 utgör Kosterbåten 1.
 2
 Del om 1 627 m av Bilen 4 utgör Osterbåten 2.
 Del om 1 784 m² av Bilen 4 utgör Koggen 1.
 Del om 1 666 m² av Bilen 4 utgör Koggen 2.
 Areal : 225 545 m²
 Avst reg 16/3-2011: Se akt ang delar om ca 2060 m² resp 11 m² av fastigheten Bilen 7 i lagf skåp
 Del om 759 m² av Bilen 4 utgör Bilen 13. A 515 och F 11
 Areal: 224 786 m² 2 0 0 9

Output file 74: TOCR – Image d

Fastigheten Oxie 1:5
 Frgl reg 19/4-79: se akt ang fastigh Oxie 1:1 m fl
 Del om 10 801,7 d av samfälligheten Oxie s:2 har tillagts fastigheten Oxie 1:5 (F
 A61, F60
 Enl. avtal -7g7g har del av fastigheten Oxie 1:5 avyttrats,

se akt ang, ca 78.000 J av fastigheten Oxie 21;2. Delen har vid frgl. reg. 15/1-81 visat sig innehålla 7.287,9 d.
Areal: 2 227 292,6
Enl. avtal A 20/1961 har del om 52.615,9 av Oxie 1:5, avs.
att bilda fastigheten Fornl ämningen 1 avyttrats. 7
Avst reg 17/3-88: se akt ang Oxie 1:1 m fl
Del om 14 225 m² a oxie 1:5 utgör Fornlämningen 2
Areal: 2 180 452 m

Output file 75: *TOCR – Image e*

Fastigheten Oxie 1:5
Bildad genom sammanläggning av fastigheterna Oxie 1:1,
Oxie 17:2, Oxie 19:1, Oxie 20:8, Oxie 25:9, Oxie 28:1 och
Oxie 29:1.
Areal efter sammanläggning: 2.231.222 d
Särskilt namn en1. reg. beslut den 15/4-92: Enedriksherg
Frgl. reg.1?/IO-.75 : se akt ang, stg 2440 m. fl. i Fosie
Delar om 8530 J av samfällda området add, 1.956,0 l av sam-
fällde området Oxie s:1, 33.507,2 av stedsägan.2834 och
22.12?,4 S av stadsägen 2467,2468 har tillagts fastigheten
Ori e 1: 5.

Output file 76: *TOCR – Image f*

Forts från
Fastigheten Oxie 1:5 4.
Del om 1 261 m av fastigheten Oxie 1:5 utgör samfälligheten
Gånggriften S:1.
Areal 2 041 038 m²
Erligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten
Oxie 1:5 avyttrats.

Output file 77: *TOCR – Image g*

Forts från
Fastigheten Oxie 1:5 4.
Del om 1 261 m av fastigheten Oxie 1:5 utgör samfälligheten
Gånggriften S:1.
Areal 2 041 038 m²
Erligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten
Oxie 1:5 avyttrats.

Output file 78: *TOCR – Image h*

i~otis ffän
Fastigheten 0#ie 1:5 4.
Del om 1 261 m av fastigheten Oxie 1:5 utgör samf811igheten
Gånggritten S:1.
nìgt avtal A 329/2006 har del om ca 23 300 m² av fastigheten
Oxie 1:5 avyttrats

Output file 79: *TOCR – Image i*

Forts från
 Fastigheten Oxie 1:5 4.
 Del om 1 261 m av fastigheten Oxie 1:5 utgör samfälligheten
 Gånggriften S:1.
 Areal 2 041 038 m²
 Eriligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten
 Oxie 1:5 avyttrats.

Output file 80: *TOCR – Image j*

Image (shornn in Appendix A)
 a b c d e
 Tesseract 16.97% 35.59% 16.72% 36.78% 16.4%
 \$ Ocrad 28.96% 57.93% 36.2% 75.78% 51.51%
 OpenOCR. 23.43% 49.01% 17.33% 7.442% 17.68%
 GOOCR. 62.62% 41.04% 77.68% 62.86% 78.53%
 OCR.opus 32.87% 15.79% 37.64% 41.52% 39.38%
 TOCR. 4.74% 4.413% 13.22% 7.207% 11.26%
 Abbyy 9.513% 29.16% 24.56% 27.03% 15.31%

Table 1: Results of the OCR scans.

80

60

40

20

0

b c d e

Image (shornn in Appendix A)
 Tesseract Ocrad OpenOCR. GOOCR.
 OCR.opus TOCR. Abbyy

Figure 2: Illustration of the scan results.

The mean error value in increasing order for each of the OCR. tools is shown below:

TOCR 8.169%
 Abbyy 21.11%
 OpenOCR 22.98%
 Tesseract 24.49%
 OCROpus 33.44%
 Ocrad 50.07%
 GOOCR 64.55%
 6

Output file 81: *TOCR – Image k*

Testing OCR tools
 Angelica Gabasio
 1 English
 This is a page with a picture to test the OCR tools.
 The page contains text in English and Swedish.
 Figure 1: This is a cat.
 2 Swedish
 Och så lite svenska tecken. Båten åker i sjön.
 Det är blött.
 1

Output file 82: *TOCR – Image l*

p cd har c acc nho
 ANGEL A ABA5 c

Output file 83: *TOCR – Image m*

The algorithm uses a database of letters and numerals to match the characters. The algorithm will calculate a matching value for each character and selects the option with the best value.[1] A way of improving the output is to check if the scanned word exists in a dictionary. If it doesn't, it is likely that some of the analyzed letters are wrong, and some of the letters may be changed so the word matches a word in the dictionary. This is done by adding a penalty value to the matching value if the analyzed character result in a word not in the dictionary. This makes it possible for another letter to get a better matching value if the word with the new letter is part of the dictionary. The letters are not always changed, sometimes a non-dictionary word gets a better matching score, even with the penalty value, than the one with the changed letters that is in the dictionary.

Even if using a dictionary may help the output accuracy it is no guarantee that the correct word will be found.[1] An example of the use of a dictionary

can be seen in Figure 1.

Other ways to get a better output is to use machine learning, where the OCR tool can be trained to recognize more languages or fonts, or to use the knowledge of which character combinations are more likely to appear. The better the image is, the more accurate the result gets. If the OCR algorithm is applied on a horizontal text with very little or no noise and the letters are well separated from each other it is possible to get a very good result. But if the image is blurry, warped or noisy, or if the characters are merged, it is likely that the output gets worse than it would otherwise.[1, 2]

The output accuracy will also depend on the image resolution. OCR software is generally better at recognizing machine writing than handwriting. In machine writing the same character always looks the same, given that the same font is used. With handwriting the characters differ a lot depending on who is writing, and even if the same person has written a text, the same character will be a little different each time it is written. Machine writing often consists of more distinct characters, and since handwriting is more varied than machine writing, it is harder for the algorithm to match the characters to its database.

3

Output file 84: *TOCR – Image n*

B.7 Abbyy

Forts (, !"X 6'
Fastigheten; Bilen 4 N
Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
Del om 836 m av fastigheten Bilen 4 utgör nybildning av
samfäll igheten Stapel bädden S:4 (se överenskommelse i avyttrings -
2
akt ang del om ca 3 020 m av fastigheten Bilen 4)
Areal: 339 546 m²
Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737 m² av fastigheten Bilen 4 utgör Galeasen 1.
2
Del om 2 229 m av fastigheten Bilen 4 utgör Galeasen 2~
Areal: 336 580 m²
V G V!
Forts från fastigheten bia ,
Bilen 4
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 m av samfälligheten Flaggskapparen S:1 har tillagts
fastigheten Bilen 4.
2
Del om 9 m av fastigheten Bilen 4 har tillagts samfäll igheten
Flaggskepparen s:1.
Areal: 332 894 m²
Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: 236 416 m²
V G V!

Output file 85: *Abbyy – Image a*

Forts (,
Fastigheten Bilen 4
Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
2
Del om 836 m av fastigheten Bilen 4 utgör nybildning av
samfäll igheten Stapel bädden S:4 (se överenskommelse i avyttrings -
2
akt ang del om ca 3 020 m av fastigheten Bilen 4)
Areal: 339 546 m²
Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737 m² av fastigheten Bilen 4 utgör Galeasen 1.
2
Del om 2 229 m av fastigheten Bilen 4 utgör Galeasen 2~
Areal: 336 580 m²
V G V!
Forts från fastigheten bia 1
Bilen 4
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
2
Del om 63 m av samfälligheten Flaggskapparen S:1 har tillagts
fastigheten Bilen 4.
2
Del om 9 m av fastigheten Bilen 4 har tillagts samfäll igheten
Flaggskepparen s:1.
Areal: 332 894 m²
Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om 96 478 m av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: 236 416 m²

V G V!

Output file 86: *Abbyy – Image b*

Fastigheten Bilen 4
Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
2
Del om 836 m² av fastigheten Bilen 4 utgör nybildning av
samfälligheten Stapel bidden S:4 (se överenskommelse i avyttrings -
2
akt ang del om ca 3 020 m² av fastigheten Bilen 4)
Areal: 339 546 m²
/st reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
il om 737 m² av fastigheten Bilen 4 utgör Galeasen 1.
2
:1 om 2 229 nr av fastigheten Bilen 4 utgör Galeasen 2~
-eal: 336 580 m²
V G V!
Forts från fastigheten Sld 7
Bilen 4
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
2
Del om 63 m² av samfälligheten Flaggskepparen S:1 har tillagts
fastigheten Bilen 4.
2
Del om 9 m av fastigheten Bilen 4 har tillagts samfälligheten
Fl aggskepparen s:1.
Areal: 332 894 m²
vst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: 236 416 m²
V G V!

Output file 87: *Abbyy – Image c*

Enligt avtal A 515 och F 11 har: del om ca 760 m² av fastigheten
2009
Bilen 4 avyttrats. Se akt ang delar om ca 2 060 m² resp 11 m²
av fastigheten Bilen 7 i lagf skåp. Avst reg 16/3-2011
Enligt avtal A 34/2010 har del om ca 2 400 m² av fastigheten
Bilen 4 avyttrats. Avst reg 15/7-10
Avst reg 15/7-10: Se avyttringsakt ang del om ca 2 400 m² av
Bilen 4.
Del om 2 428 m² av Bilen 4 utgör fastigheten Klyvaren 1. A 34/201
Del om 1 312 m² av Bilen 4 utgör fastigheten Klyvaren 2.
Areal: 332 840 m² ----- Forts sid 7
Avst reg 25/1-2011: Se akt ang Bilen 4 i lagf skåp
! Del om 1 017 m² av Bilen 4 utgör Klippen 1.
I Del om 813 m² av Bilen 4 utgör Klippen 2.
[Del om 631 m² av Bilen 4 utgör Klippen 3.
i Del om 1 078 m² av Bilen 4 utgör Klippen 4.
Del om 2 255 m² av Bilen 4 utgör ;<osterbåten 1.
i 2
Del om 1 627 m² av Bilen 4 utgör JÖsterbåten 2.
f Del om 1 784 m² av Bilen 4 utgör Koggen 1.
I Del om 1 666 m² av Bilen 4 utgör Koggen 2.
Areal: 225 545 m²
I
Avst reg 16/3-2011: Se akt ang delar om ca 2060 m² resp 11 m² av

fastigheten Bilen 7 i lagf skåp
De? om 759 m² av Bilen 4 utgör Bilen 13. A 515 och F 11
Areal: 224 786 m² 2 0 0 9

Output file 88: *Abbyy – Image d*

Fastigheten Oxie 1:5
Frgl reg 19/4-79: se akt ang fastigh Oxie 1:1 a fl
Del om 10 801,7 rfi av samfälligheten Oxie s:2 har tillagts
fastigheten Oxie 1:5 (| f^/j^)
, . A61, F60
Enl. avtal har del av fastigheten Oxie 1:5 avyttrats,
se akt ang. ca 78.000 nf av fastigheten Oxie 21; 2. Delen har
vid frgl. reg. 15/1-81 visat sig innehålla 7.287,9 ma.
Areal: 2 22 7 2 92,6 af
Enl. avtal A 20/1981 har del om 3 2,615,9 m² av Oxie 1:5, avs.
att bilda fastigheten Porn1 amningen 1 avyttrats. %-*s)
Avst reg 17/3-88: se akt ang Oxie 1:1 m fl
Del om 1^ 225 m² av Oxie 1:5 utgör Fornlämningen 2
Areal: 2 180 ^52 m²

Output file 89: *Abbyy – Image e*

57, c
Fastigheten Oxie 1:5
Bildad genom sammanläggning av fastigheterna Oxie 1:1,
Oxie 17:2, Oxie 19:1, Oxie 20:8, Oxie 25:9, Oxie 28:1 och
Oxie 29:1.
.Areal efter sammanläggning: 2.231.222 nf
Särskilt namn enl. reg. beslut den 15/4-92: i I ksb.erc}
Frgl. reg.17/10-75 : se akt ang. stg 2440 m. fl. i Fosie
Delar om 85,0 m av samfällda området adj., 1.956,0 nf av sam-
fällda området Oxie s:1, 33.507,2 nf av stadsägan. 2834 och
22.122,4 nf ?.v stadsägan 2^67, 2468 har tillagts fastigheten
Oxie 1:5.

Output file 90: *Abbyy – Image f*

Forts från
Fastigheten Oxie 1:5
4,
Del om 1 261 m jiv^fastigheten Oxie 1:5 utgör samfäll i gheten
Gång^riften S:1.
Areal: 2 041 038 m²
Enligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten
Oxie 1 :5 avyttrats._____._____\\

Output file 91: *Abbyy – Image g*

Forts från
Fastigheten Oxie 1:5
4,
Del om 1 261 m jiv^fastigheten Oxie 1:5 utgör samfäll i gheten
Gång^riften S:1.
Areal: 2 041 038 m²
Enligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten

Ox i e 1 :5 avyttrats._____.__________

Output file 92: *Abbyy – Image h*

Output file 93: *Abbyy – Image i*

Forts från
Fastigheten Oxie 1:5
Del om 1 261 m av fastigheten Oxie 1:5 utgör samfälligheten
Gånggriften S:1.
Areal: 2 041 038 m²m
Enligt avtal A 329/2006 har del om ca 23 300 m av fastigheten
Oxie 1:5 avyttrats._____._____.1_ r> __

Output file 94: *Abbyy – Image j*

a	b	c	d	e	pendix A)
Tesseract	16.97%	35.59%	16.72%	36.78%	16.4%
Ocrad	28.96%	57.93%	36.2%	75.78%	51.51%
OpenOCR	23.43%	49.01%	17.33%	7.442%	17.68%
GOCR	62.62%	41.04%	77.68%	62.86%	78.53%
OCRopus	32.87%	15.79%	37.64%	41.52%	39.38%
TOCR	4.74%	4.413%	13.22%	7.207%	11.26%
Abbyy	9.513%	29.16%	24.56%	27.03%	15.31%

Table 1: Results of the OCR seans.

Til

1

b c d

Image (shown in Appendix A)

ils is shown

III Tesseract 11 Ocrad III OpenOCR GOCR

11 OCRopus 11 TOCR Abbyy

Figure 2: Illustration of the sean results.

TOCR 8.169%

Abbyy 21.11%

OpenOCR 22.98%

Tesseract 24.49%

OCRopus 33.44%

Ocrad 50.07%

GOCR 64.55%

f)

Output file 95: *Abbyy – Image k*

Testing OCR tools

Angelica Gabasio

1 English

This is a page with a picture to test the OCR tools.

The page contains text in English and Swedish.

W

II

Figure 1: This is a cat.

2 Swedish

Och så lite svenska, tecken. Båten åker i sjön.

Det är blött.
i

Output file 96: *Abbyy – Image l*

optical CWarae^e-r rtcoanvilon
2013

Output file 97: *Abbyy – Image m*

The algorithm uses a database of letters and numerals to match the characters. The algorithm will calculate a matching value for each character and selects the option with the best value. [1]
A way of improving the output is to check if the scanned word exists in a dictionary. If it doesn't, it is likely that some of the analyzed letters are

wrong, and some of the letters may be changed so the word matches a word in the dictionary. This is done by adding a penalty value to the matching value if the analyzed character result in a word not in the dictionary. This makes it possible for another letter to get a better matching value if the word

with
the new letter is part of the dictionary. The letters are not always changed
,

sometimes a non-dictionary word gets a better matching score, even with the penalty value, than the one with the changed letters that is in the dictionary.

Even if using a dictionary may help the output accuracy it is no guarantee that the correct word will be found. [1] An example of the use of a dictionary

can be seen in Figure 1.

Other ways to get a better output is to use machine learning, where the OCR tool can be trained to recognize more languages or fonts, or to use the knowledge of which character combinations are more likely to appear.

The better the image is, the more accurate the result gets. If the OCR algorithm is applied on a horizontal text with very little or no noise and the

letters are well separated from each other it is possible to get a very good result. But if the image is blurry, warped or noisy, or if the characters are

merged, it is likely that the output gets worse than it would otherwise.[1,
2]

The output accuracy will also depend on the image resolution.

OCR software is generally better at recognizing machine writing than handwriting. In machine writing the same character always looks the same, given that the same font is used. With handwriting the characters differ a lot

depending on who is writing, and even if the same person has written a text, the same character will be a little different each time it is written.

Machine

writing often consist of more distinct characters, and since handwriting is more varied than machine writing, it is harder for the algorithm to match the characters to its database.

Output file 98: *Abbyy – Image n*

B.8 Leadtools

64
Forts i
Fastigheter Bilen 4
Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
Del om 836 m² av fastigheten Bilen 4 utgör nybildning av
1
samfälligheten Stapelbädden S:4 (se överenskommelse i avyttrings-,
akt ang del om ca 3 020 m² av fastigheten Bilen 4)
T Areal: 339 546 m²
Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737m² av fastigheten Bilen 4 utgör Galeasen 1.
i 1 Del om 2 229 m² av fastigheten Bilen 4 utgör Galeasen
Areal: 336 580 m²
. V G V!
Förts från fastigheten Sid 7
Bilen 4
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 m² av samfälligheten Flaggskiepparen S:1 har tillagts
fastigheten Bilen 4.
Del om 9 m² av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s: 1.
Areal: 332 894 m²
Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: 236 416 m²
V G V!

Output file 99: *Leadtools – Image a*

6 ,
Forts ^
Fastigheten Bilen 4
Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
Del om 836 m² av fastigheten Bilen 4 utgör nybildning av
I
samfälligheten Stapelbädden 5:4 (se överenskommelse i avyttrings-,
! 1
akt ang del om ca 3 020 m² av fastigheten Bilen 4)
Areal: 339 546 m²
Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737m² av fastigheten Bilen 4 utgör Galeasen le
i 1 Del om 2 229 m² av fastigheten Bilen 4 utgör Galeasen 2.
Areal: 336 580 m²
V G V!
Förts från fastigheten Sid i
Bilen 4
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 m² av samfälligheten Flaggskiepparen 5:1 har tillagts
fastigheten Bilen 4.
Del om 9 m² av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s: 1.
Areal: 332 894 m²
Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: 236 416 m²
V G V!

Output file 100: *Leadtools – Image b*

Forts
 Fastigheten Bilen 4
 64
 fulseg 19/8-09: Se akt ang Bilen 4 i Lagt' skåp
 Del om 836 m² av fastigheten Bilen 4 utgör nybildning av
 ,
 samfälligheten Stapelbädden S:4 (se överenskommelsei avyttrings- ,
 i
 i akt ang del om ca 3 020 m² av fastigheten Bilen 4)
 Areal: 339 546 m²
 '52..s1,22,2219-09: Se akt ang Bilen 4 i lagf skåp
 /
 Del om 737 m² av fastigheten Bilen 4 utgör Galeasen 1.
 Del om 2 229 m² av fastigheten Bilen 4 utgör Galeasen 2.,
 Areal: 336 580 m²
 V G V!
 Forts frän fastigheten Sid 7
 Bilen 4
 Fr i re 30/7-10: Se akt ang Bilen 4 i lagt' skåp
 Del om 63 m² av samfälligheten Flaggskepparen 5:1 har tillagts
 fastigheten Bilen 4.
 D i el om 9 m² av fastigheten Bilen 4 har tillagts samfälligheten
 J
 nagskepparen 5:1.
 Areal: 332 894 m²
 Avst reg 23/11-2010: Se akt ang Bilen 4 i hyr skåp
 Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten
 Bilen 12.
 Areal: 236 416 m²
 V G V!

Output file 101: Leadtools - Image c

Enligt avtal A 515 och F 11 har del om ca 760 m² av fastigheten
 2009
 Bilen 4 avyttrats. Se akt ang delar omca 2 060} resp 11 m² –
 av fastigheten Bilen 7 i lagf skåp. Avst reg 16/3-2011
 Enligt avtal A 34/2010 har del om ca 2 400 m² av fastigheten
 I
 Bilen 4 avyttrats. Avst reg 15/7710
 Avst reg 15/7-10: Se avyttringsakt ang del om ca 2 400 m² av
 Bilen 40
 1 Del om? 428 m² av Bilen 4 utgör fastigheten Klyvaren 1. A 34/201
 -- I Del om 1 312 m² av Bilen 4 utgör fastigheten Klyvaren 24
 t
 Areal: 33? 840 m²
 Fortssid7
 Avst reg 25/1-2011: Se akt ang Bilen 4 i lagf skåp
 Del om 1 017 m² av Bilen 4 utgör Klippern 1.
 Del om 813 m² av Bilen 4 utgör Klippern 2.
 Del om 631 m² av Bilen 4 utgör Klippern 3.
 Del om 1 078 m² av Bilen 4 utgör Klippern 4.
 Del om 2 255 m² av Bilen 4 utgör Kosterbåten 1.
 2
 Del om 1 627 m² av Bilen 4 utgör K)sterbåten 2.
 Del om 1 784 m² av Bilen 4 utgör Koggen 1.
 Del om 1 666 m² av Bilen 4 utgör Koggen 2.
 Areal: 225 545 m²
 Avst reg 16/3-2011: Se akt ang delar om ca 2060 m² resp 11 m² av
 fastigheten Bilen 7 i lagf skåp
 1 Del om 759 m² av Bilen 4 utgör Bilen 13. A 515 och F 11

Areal: 224 786m²
2 0 09

Output file 102: *Leadtools – Image d*

Fastigheten Oxie 1:5
Frgl reg 19/4-79: se akt ang fastigh Oxie 1:1 m fl
av samfälligheten Oxie s:2 har tillagts
Del om 10 801 97
d
fastigheten Oxie 1:5
(y Wiffl)
590r ,o,cpc
A61, F60
Enl. avtal har del av fastigheten Oxie 1:5 avyttrats,
1979
se akt ang. ca 78.000 å av fastigheten Oxie 21;2. Delen har
vid frgl. reg. 15/1-81 visat sig innehålla 7.2879 1112.
Areal: 2 227 292,6 ni
Enl. avtal A 20/1981 har del om 32.615,9 å av Oxie 1:5, avs.
att bilda fastigheten Fornlämningen 1 avyttrats.6' 7 '21 c4'sv)
Avst reg 17/3-88: se akt ang Oxie 1:1 m fl
Del om 14 225 m 2 ax Oxie 1:5 utgör Fornlämningen 2
452 m'
2 180
Areal:

Output file 103: *Leadtools – Image e*

4
Fastigheten Oxie 1:5
Bildad genom sammanläggning av fastigheterna Oxie 1:1,
Oxie 17:2, Oxie 19:1, Oxie 20:8, Oxie 25:9, Oxie 28:1 och
Oxie
,AL25.2.22122=haailW 2.231.222 af
Särskilt namn en1 teg. beslut den 15/4-92: Uzedriksherg
a : m 6 s 1Fx. e.1710--.12 a.: se akt ang. stg 2440 m.fl. i Fosie
av samfälidg området adj, 1.956,0 Å av sams
Delar om 85,0
J
fälld området Oxie s:1 ? 55.507,2 rlf av stadsägan.2854 och
92.122,4 d sv stadsägan 2467,2468 har tillagts fastigheten
Oxie 1:5.

Output file 104: *Leadtools – Image f*

Forts från
Fastigheten Oxie i:5 40
Del om 1 261
ffi
av fastigheten Oxie 1:5 utgör samfälligheten
Gånggriften S:1.
Areal: 2 041 038 m²
1 Enligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten
Oxie 1:5 avyttrats. ___.
_e__

Output file 105: *Leadtools – Image g*

Forts från
Fastigheten Oxie 1:5 4.
Del omr 261
ffl av fastigheten Oxie
Gånggriften S:1.
Areal: 2 041 038 m²
1 Enligt avtal A 329/20
Oxie 1:5 avyttrats. 1.

Output file 106: *Leadtools – Image h*

Output file 107: *Leadtools – Image i*

Forts frän
Fastigheten Oxie 1:5 4.
I Del omr 261 m² av fastigheten Oxie 1:5 utgör samfälligheten
Gånggri ften S:1.
Areal: 2 041 038 m²

'Ungt avtal A 329/2006 har del om ca 23 30 m² av fo, - igheten .
0 4
Oxie 1:5 avyttrats. An, .
?,
8 1
. " 9
1
1
- - - -
e
. ,
-
Ta,
/1 e
a1 ,
O

Output file 108: *Leadtools – Image j*

Image (shown in Appendix A)
a

```

Tesseract 16.97% 35.59% 16.72% 36.78% 16.4%
Ocrad 28.96% 57.93% 36.2% 75.78% 51.51%
OpenOCR 23.43% 49.01% 17.33% 7.442% 17.68%
GOCR 62.62% 41.04% 77.68% 62.86% 78.53%
OCRopus 32.87% 15.79% 37.64% 41.52% 39.38%
TOCR 4.74% 4.413% 13.22% 7.207% 11.26%
Abbyy 9.513% 29.16% 24.56% 27.03% 15.31%

```

Table 1: Results of the OCR scans.

```

80
f1
(1) 60
z:s.) 1111! I
(1)
(1) gt. 40 I 11 1 I
LI
20-
I
0
a
```

```

Image (shown in Appendix A)
I I Tesseract II Ocrad I I OpenOCR GOCR
II OCRopus II TOCR I I Abbyy
```

Figure 2: Illustration of the scan results.

The mean error value in increasing order for each of the OCR tools is shown below:

```

TOCR 8.169%
Abbyy 21.11%
OpenOCR 22.98%
Tesseract 24.49%
OCRopus 33.44%
Ocrad 50.07%
GOCR 64.55%
6
```

Output file 109: *Leadtools – Image k*

```

Testing OCR tools
Angelica Gabasio
1 English
This is a page with a picture to test the OCR tools.
The page contains text in English and Swedish.
il , fra
: \
i , 4
"niors
Eigure 1: This is a cat.
2 Swedish
Och så lite svenska tecken. Båten åker i sjön.
Det är blött.
1
```

Output file 110: *Leadtools – Image l*

```

a P k. i cak- C bara:åar. re_co_n_g 'Uön
L. ;
ANGEL\ CA GIABN5 0
2.013
```

Output file 111: *Leadtools – Image m*

The algorithm uses a database of letters and numerals to match the characters. The algorithm will calculate a matching value for each character and selects the option with the best value. [1]

A way of improving the output is to check if the scanned word exists in a dictionary. If it doesn't, it is likely that some of the analyzed letters are

wrong, and some of the letters may be changed so the word matches a word in the dictionary. This is done by adding a penalty value to the matching value if the analyzed character result in a word not in the dictionary. This makes it possible for another letter to get a better matching value if the word with

the new letter is part of the dictionary. The letters are not always changed

sometimes a non -dictionary word gets a better matching score, even with the penalty value, than the one with the changed letters that is in the dictionary.

Even if using a dictionary may help the output accuracy it is no guarantee that the correct word will be found. [1] An example of the use of a dictionary

can be seen in Figure 1.

Other ways to get a better output is to use machine learning, where the OCR tool can be trained to recognize more languages or fonts, or to use the knowledge of which character combinations are more likely to appear.

The better the image is, the more accurate the result gets. If the OCR algorithm is applied on a horizontal text with very little or no noise and the

letters are well separated from each other it is possible to get a very good result. But if the image is blurry, warped or noisy, or if the characters are

merged, it is likely that the output gets worse than it would otherwise.[1, 2]

The output accuracy will also depend on the image resolution.

OCR software is generally better at recognizing machine writing than handwriting. In machine writing the same character always looks the same, given that the same font is used. With handwriting the characters differ a lot

depending on who is writing, and even if the same person has written a text, the same character will be a little different each time it is written.

Machine

writing often consist of more distinct characters, and since handwriting is more varied than machine writing, it is harder for the algorithm to match the characters to its database.

3

Output file 112: *Leadtools – Image n*

B.9 OCR API Service

Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
2
Del om 836 m av fastigheten Bilen 4 utgör nybildning av
samfälligheten Stapel bädden S:4 (se överenskommelse i avyttrings -
2
akt ang del om ca 3 020 m av fastigheten Bilen 4)
Areal: 339 546 m2
Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737 m2 ay fastigheten Bilen 4 utgör Galeasen 1.
Del om 2 229 m2 av fastigheten Bilen 4 utgör Galeasen 2.
Areal: 336 580 m2
V G V!
Forts från fastigheten 51d 7
Bilen 4
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
2
Del om 63 m av samfälligheten Flaggskiepparen S:1 har tillagts
fastigheten Bilen 4.
2
Del om 9 m av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s:1.
Areal: 332 894 m2
Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om 96 478 m2 av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: 236 416 m2
V G V!

Output file 113: *OCR API Service – Image a*

Forts s
Fastigheterf Bilen 4 \)
6,
Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skap
2
Del om 836 m av fastigheten Bilen 4 utgör nybildning av
samfälligheten Stapel bädden S:4 (se överenskommelse i avyttrings -
2
akt ang del om ca 3 020 m av fastigheten Bilen 4)
Areal: 339 546 m2
Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737 rr 2 av fastigheten Bilen 4 utgör Galeasen 1.
Del om 2 229 m2 av fastigheten Bilen 4 utgör Galeasen 2
Areal: 336 580 m2
V G V!
Forts från fastigheten 2)1 a 1
Bilen 4
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 m av samfälligheten Flaggskiepparen S:1 har tillagts
fastigheten Bilen 4.
2
Del om 9 m av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s:1.
Areal: 332 894 m2
Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
o
Del om 96 478 m av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: 236 416 m2

V G V!

Output file 114: *OCR API Service – Image b*

Forts
Fastigheten Bilen 4
6,
Frgl reg 19/8-09: Se akt ang Bilen 4 i lagf skåp
2
Del om 836 m av fastigheten Bilen 4 utgör nybildning av
samfälligheten Stapel bädchen S:4 (se överenskommelse i av
2
akt ang del om ca 3 020 m av fastigheten Bilen 4)
Areal: 339 546 m²
Avst reg 30/9-09: Se akt ang Bilen 4 i lagf skåp
Del om 737 m² ay fastigheten Bilen 4 utgör Galeasen 1.
Del om 2 229 m² av fastigheten Bilen 4 utgör Galeasen 2
Areal: 336 580 m²
V G V!
Forts från fastigheten Sld 7
Bilen 4
Frgl reg 30/7-10: Se akt ang Bilen 4 i lagf skåp
Del om 63 m av samfälligheten Flaggskepparen S:1 har tillagts
fastigheten Bilen 4.
2
Del om 9 m av fastigheten Bilen 4 har tillagts samfälligheten
Flaggskepparen s:1.
Areal: 332 894 m²
t
Avst reg 23/11-2010: Se akt ang Bilen 4 i lagf skåp
Del om 96 478 m² av fastigheten Bilen 4 utgör fastigheten
Bilen 12.
Areal: 236 416 m²
V G V!

Output file 115: *OCR API Service – Image c*

Enligt avtal A 515 och F 11 har del om ca 760 nr av fastigheten
2009 "
Bilen 4 avyttrats. Se akt ang delar om ca 2 060 m² resp H m²
av fastigheten Bilen 7 i lagf skåp. Avst reg 16/3-2011
Enligt avtal A 34/2010 har del om ca 2 400 nr av fastigheten
Bilen 4 avyttrats. Avst reg 15/7-10
Avst reg 15/7-10: Se avyttringsakt ang del om ca 2 400 m² av
Bilen 4.
Del om 2 428 m² av Bilen 4 utgör fastigheten Klyvaren 1. A 34/201
Del om 1 312 m² av Bilen 4 utgör fastigheten Klyvaren 2.
Areal: 332 840 m²
----- Forts
sid 7
Avst reg 25/1-2011: Se akt ang Bilen 4 i lagf skåp
I Del om 1 017 m av Bilen 4 utgör Klippan 1.
I 2
I Del om 813 m av Bilen 4 utgör Klippan 2.
! Del om 631 m av Bilen 4 utgör Klippan 3.
i Del om 1 078 m² av Bilen 4 utgör Klippan 4.
f 2
i Del om 2 255 m av Bilen 4 utgör Kosterbåten 1.
I Del om 1 627 m av Bilen 4 utgör Kosterbåten 2.
Del om 1 784 m² av Bilen 4 utgör Koggen 1.
Del om 1 666 m² av Bilen 4 utgör Koggen 2.

Areal: 225 545 m²
Avst reg 16/3-2011: Se akt ang delar om ca 2060 m² resp 11 m² av fastigheten Bilen 7 i lagf skåp
I De? om 759 m² av Bilen 4 utgör Bilen 13. A 515 och F 11
Areal: 224 786 m² 0 n n a

Output file 116: *OCR API Service – Image d*

Fastigheten Oxie 1:5
Frgl reg 19/4-79: se akt ang fastigh Oxie 1:1 m f1
Del om 10 801,7 nf av samfälligheten Oxie s:2 har tillagts fastigheten Oxie 1:5 (j)
^ ^ 3^ S~#Ö. S- sr.
Enl. avtal ^^g^g^ ^ar öel av fastigheten Oxie 1:5 avyttrats, se akt ang. ca 78.000 nf av fastigheten Oxie 21; 2. Delen har vid frgl. reg. 15/1-81 visat sig innehålla 7.287,9 m.
Areal: 2 227 292,6 rf
Enl. avtal A 2Q/1981 har del om 3 2,615,9 m* av Oxie 1:5, avs. att bilda fastigheten Pornl amningen 1 avyttrats.
Avst reg 17/3-88: se akt ang Oxie 1:1 m fl
Del om Ik 225 2 av Oxie 1:5 utgör Fornlämningen 2
: 2 180 k\$2 m2

Output file 117: *OCR API Service – Image e*

Fastigheten Oxie 1:5
Bildad genom sammanläggning av fastigheterna Oxie 1:1, Oxie 17:2, Oxie 19:1, Oxie 20:8, Oxie 25:9, Oxie 28:1 och Oxie 29:1.
Areal efter sammanläggning: 2.231.222 m
Särskilt namn enl. reg. beslut den 15/4-92:.. Fxe.driksb.erg^
Frgl. reg, 17/10-75 : se akt ang. stg 2440 m. fl. i Fosie
Delar cm 85,0 af av samfällda området adj., 1,956,0 ra2 av samfällda området Oxie s:1, 33.507,2 i av stadsägan. 2834 och 22.122,4 i av stadsägan 2467,2468 har tillagts fastigheten Oxie 1:5.

Output file 118: *OCR API Service – Image f*

Forts från
Fastigheten Oxie 1:5
Del qm l_261 m av fastigheten
Gång^riften S:1.
Areal: 2 041 038 m
Enl igt avtal A 329/2006 har del
Oxie 1:5 avyttrats._____
4.
3xiel:5 utgör samfäll igheten
2
om ca 23 300 m av fastigheten

Output file 119: *OCR API Service – Image g*

Forts från
Fastigheten Oxie 1:5
4.

Del om 1 261 m av fastigheten Oxie 1:5 utgör samfäl ligheten
 Gång^riften S:1.
 Areal: 2 041 038 m²
 Enligt avtal A 329/2006 har del om ca 23 300 m² av fastigheten
 Oxie 1 ^y^yttrats^j

Output file 120: *OCR API Service – Image h*

Output file 121: *OCR API Service – Image i*

Forts från
 Fastigheten Oxie 1:5
 Del om 1 261 m av fastigheten Oxie 1:5 jJtgör samfälligheten
 Gånggriften S:1
 Areal: 2 041 038 m²
 2
 4
 1 Enligt avtal A 329/2006 har del om ca 23 300 m² av f
 Oxie 1:5 avvttrats.

Output file 122: *OCR API Service – Image j*

Image	(shown in Appendix A)	a	b	c	d	e	
Tesseract		16.97%	35.59%	16.72%	36.78%	16.4%	
Ocrad	28.96%	57.93%	36.2%	75.78%	51.51%		
OpenOCR	23.43%	49.01%	17.33%	7.442%	17.68%		
GOCR	62.62%	41.04%	77.68%	62.86%	78.53%		
OCRopus	32.87%	15.79%	37.64%	41.52%	39.38%		
TOCR	4.74%	4.413%	13.22%	7.207%	11.26%		
Abbyy	9.513%	29.16%	24.56%	27.03%	15.31%		

Table 1: Results of the OCR seans.

I1 Tesseract I1 Ocrad I1 OpenOCR GOCR

I1 OCRopus II TOCR Abbyy

Figure 2: Illustration of the scan results.

The mean error value in increasing order for each of the OCR tools is shown below:

TOCR	8.169%
Abbyy	21.11%
OpenOCR	22.98%
Tesseract	24.49%
OCRopus	33.44%
Ocrad	50.07%
GOCR	64.55%

6

Output file 123: *OCR API Service – Image k*

Testing OCR tools
 Angelica Gabasio
 1 English
 This is a page with a picture to test the OCR tools.
 The page contains text in English and Swedish.
 Figure 1: This is a cat
 2 Swedish

Och så lite svenska tecken. Båten åker i sjön.
Det är blött.
1

Output file 124: *OCR API Service – Image l*

Output file 125: *OCR API Service – Image m*

The algorithm uses a database of letters and numerals to match the characters. The algorithm will calculate a matching value for each character and selects the option with the best value. [1]

A way of improving the output is to check if the scanned word exists in a dictionary. If it doesn't, it is likely that some of the analyzed letters are

wrong, and some of the letters may be changed so the word matches a word in the dictionary. This is done by adding a penalty value to the matching value if the analyzed character result in a word not in the dictionary. This makes it possible for another letter to get a better matching value if the word

with
the new letter is part of the dictionary. The letters are not always changed
,

sometimes a non-dictionary word gets a better matching score, even with the penalty value, than the one with the changed letters that is in the dictionary.

Even if using a dictionary may help the output accuracy it is no guarantee that the correct word will be found. [1] An example of the use of a dictionary

can be seen in Figure 1.

Other ways to get a better output is to use machine learning, where the OCR tool can be trained to recognize more languages or fonts, or to use the knowledge of which character combinations are more likely to appear.

The better the image is, the more accurate the result gets. If the OCR algorithm is applied on a horizontal text with very little or no noise and the

letters are well separated from each other it is possible to get a very good result. But if the image is blurry, warped or noisy, or if the characters are
merged, it is likely that the output gets worse than it would otherwise. [1,
2]

The output accuracy will also depend on the image resolution.

OCR software is generally better at recognizing machine writing than handwriting. In machine writing the same character always looks the same, given that the same font is used. With handwriting the characters differ a lot

depending on who is writing, and even if the same person has written a text, the same character will be a little different each time it is written.

Machine writing often consist of more distinct characters, and since handwriting is more varied than machine writing, it is harder for the algorithm to match the characters to its database.

3

Output file 126: *OCR API Service – Image n*