

Improving the quality of the text

**A pilot project to assess and correct the OCR in a
multilingual environment**

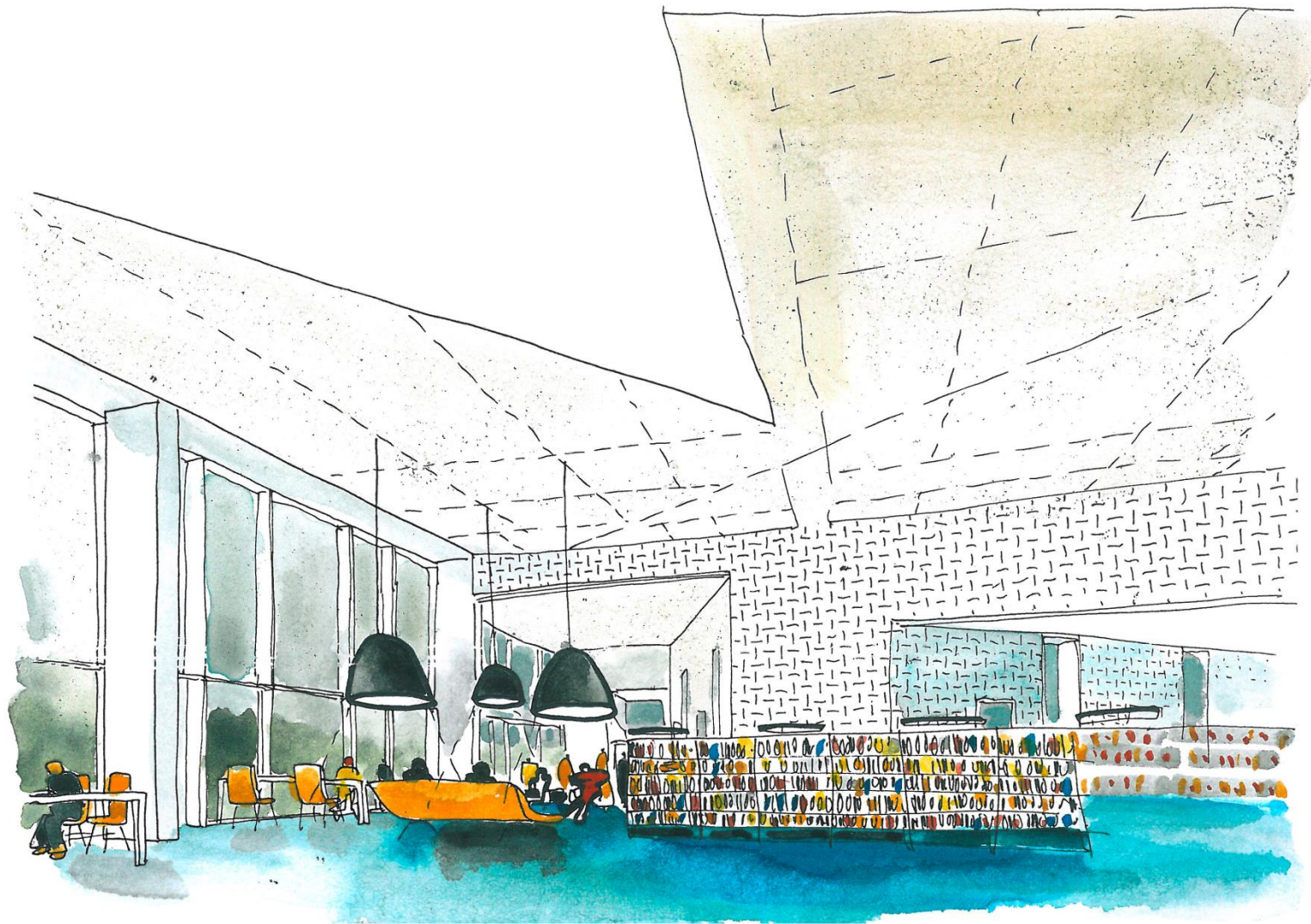
Yves Maurer

National library of Luxembourg

Outline

- The newspaper collection
- Why OCR is important
- How do you measure OCR quality?
- How to re-run OCR on METS-ALTO
- Results

Bibliothèque nationale de Luxembourg



Bibliothèque nationale de Luxembourg

- Origins in 1798
- Multiple missions: patrimonial, public and scientific
- Responsible for national library network and national consortium
- 77 FTE staff
- Digitization started in 2003 and OCR in 2006

eluxemburgensia

- 39 newspaper titles
- 612 000 pages
- 1704 – 2007
 - German
 - French
 - Luxembourgish
 - English
 - ...

+ books, posters,
postcards etc.



Find an article

★ Previous searches

➔ Advanced search



Find an issue

By Date

- 1950 ▾ +

- January ▾ +

Mo	Tu	We	Th	Fr	Sa	Su
						1
2	3	4	5	6	7	8
9	10	11	12	13	14	15
16	17	18	19	20	21	22
23	24	25	26	27	28	29
30	31					

Please note: from May 10, 1940 to September 10, 1944 the national press was controlled by the Nazi occupier.

By Newspaper



➔ d'Lëtzebuenger Land

➔ Luxemburger Wort

➔ Tageblatt

➔ Digitized books

➔ Other documents...

METS/ALTO with articles

Bibliothèque nationale de Luxembourg
Nationalbibliothek

Luxemburger Wort 1923-03-29_01

dresden

Luxemburger Wort 1923-03-29_01

Aus dem deutschen Reich.

Ausbreitungen in Dresden.

Steinlohlenpreise werden auf den 1. April um etwa 10 Prozent, die Brickettpreise um 15 Prozent herabgesetzt werden.

Ausbreitungen in Dresden.

Dresden, 28. März. (Europapress.) Im Anschluß an eine Erwerbslosenversammlung, die von anderen Erwerbslosen gesprengt wurde, ereigneten sich gestern schwere Tumulte. Die Demonstranten zogen vor das Rathaus und vor das Polizeipräsidium, wo sie durch die Polizei auseinandergejagt wurden. Darauf versuchten sie die Kaufläden in den umliegenden Straßen zu plündern, jedoch hatten die Inhaber die Läden vorher geschlossen.

Aus Frankreich.

Kammerkommission für Auswärtiges.

Paris, 28. März. (Havas.) Die Kammerkommission für auswärtige Angelegenheiten versammelte sich heute

wal
von
bish

sehe
gebe
eine
sich
stän
lief

Des
war
lom

nich
fahr
frei
han

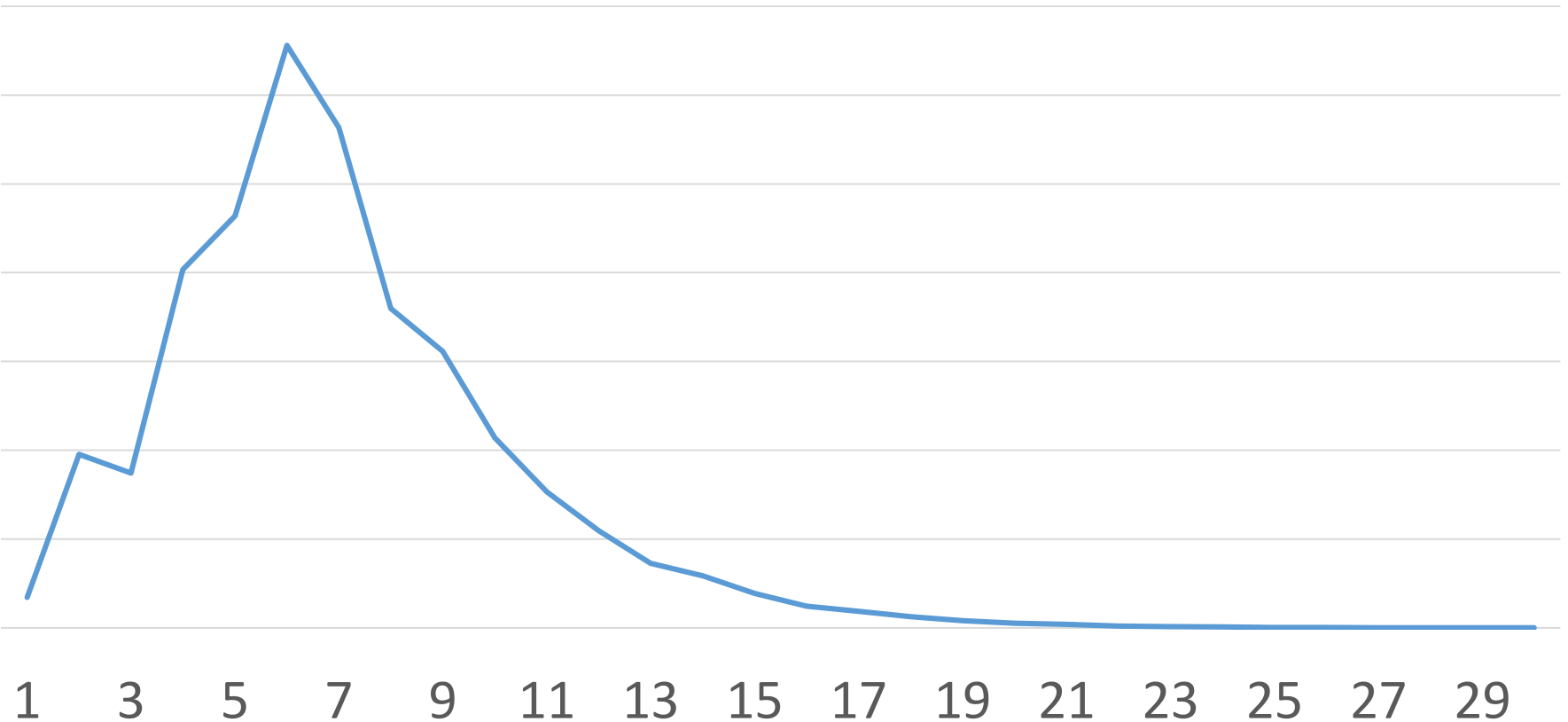
Why OCR is important

- Large collections are:
 - Impossible to browse in their entirety
 - Hard to understand for casual user

 Full text search is a useful tool

- Users expect to get all articles that match
- Text mining works when the data is "clean"
- Users search for long words (median length: 7)
- A single wrong character makes the word impossible to find

#characters in Search Terms used by Users



Data source: All searches run on www.eluxemburgensia.lu from 2010 to 2016

Measuring OCR quality

Ideally, you would compare with “correct output”

Not feasible, so use

- Word confidence
 - Measured by OCR software itself or based on character confidence
- Vocabulary growth
 - Need a large amount of text
- Dictionaries
 - For each language/variant
 - Not complete

Multilingual content

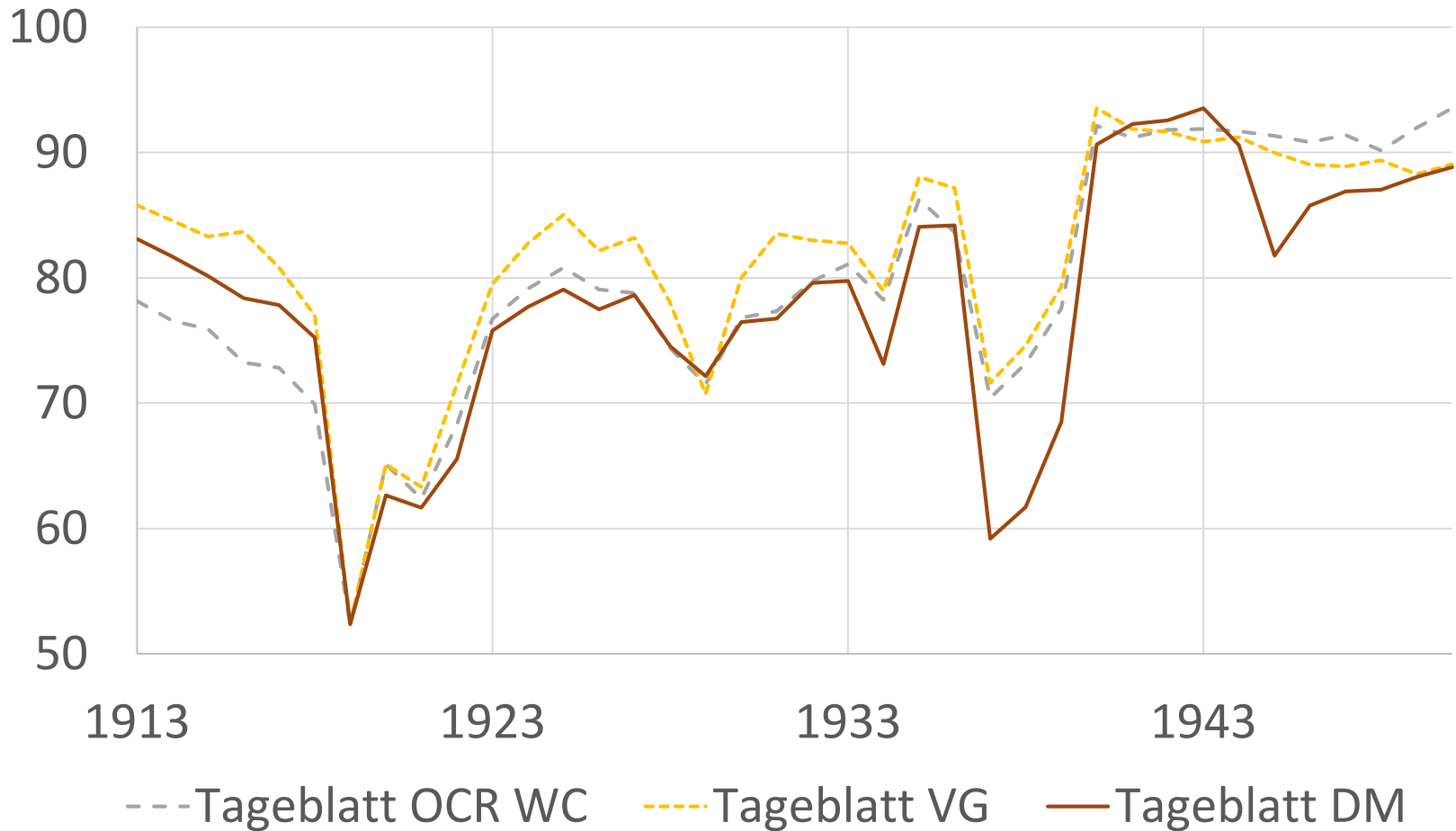
Red :
French

Blue :
Luxembourgish

Green:
English

<p>MAISON BERNHARD-SCHROEDER SUCCESSEUR VINS ET LIQUEURS Obercorn, fermée du 22 juin au 2 juillet 7294</p>	<p>Hère-Coiffeur-Salon opmachen. Em gene'gten Zo'sproch biéd den Edmond THULL.</p>	<p>traueren d'UNIO'N an de POMPIERSCORPS v. He'en. 4557</p>
<p>Achtong! Ve'hbesetzer a Baueren. No 5jähregem Exil erem dohém, délen ech de Ve'h- besetzer mat, dass ech erem ewe' fre'er mat NOTZ- an ZUCHTVE'H den Handel weiterbedreiwen. Zo' gleicher Zeit délen ech Iech mat, dass ech vum Stat, am Hollerecher Schluechthaus als Ve'hkom- missionär ernannt sin. Dir kónnst Ert fett Ve'h bei mir umellen. 4469 An der Hoffnung, daß der mir Ert Vertraue schenkt ewe' emmer, gre'ßen ech Iech mat aller Hochachtong Den Israëls Décken BO'NEWEG, Letzeburgerstroß No. 65. Tel.: Schluechthaus 58-32.</p>	<p>FOURRURES Maison de Gros - cherche clients pour le Luxembourg. Ecrire: Armand GRAULS, Bruxelles, rue Grétry 26. 4545</p>	<p>Marcel NEISELER, Zenn- techniker, gefal fir seng le'w Hémecht bei Kiel, den 3. Mè 1945, am ble'enden Alter vun 18¹/₂ Jor. Leichen- dengscht e Mittwoch, de 27. Juni, em 10 Auer, zo' Scheffleng. 7327 Famill Neiseler-Ludovicy.</p>
	<p>OUVERTURE à ESCH-Alzette</p>	<p>Sisy SCHUH, gestuerwen durch Fliegerugreff am K. H. D. Pforzheim, den 23. Fe- bruar 1945, am Alter vun 20 Jor. Feierleche Leichen- dengscht e Mèndeg, de 25. Juni, em 10.10 Auer zo' Kél, 7301 Familjen Schuh-Herzig.</p>
	<p>School of Languages Helen WIES-YULL (engl. Teacher) ENGLISH - FRENCH FRENCH and GERMAN for Allied Troops ESCH, 66, Av. de la Gare. Inscriptions pour les cours: les mercredis de 2 à 7 heures les samedis de 3 à 5 heures 4225</p>	<p>Jengy EDERT, Man vum Madeleine MICHAELY, ge- stuerwen fir seng Hémecht am KZ. Dachau, den 8. Abröl 1945, am Alter vu 55 Jor. Feierleche Leichen- dengscht e Mittwoch, de 27. Juni, em 10 Auer zo' Pe'teng. 4578</p>
	<p>RAJEUNISSEZ ! Redevenez souple et alerte comme à 20 ans, en élimi- nant l'excès d'acide urique accumulé dans votre orga-</p>	<p>Remerciements Kun HURST, anc. officier du génie aux Indes Néer- landaises. La messe de six semaines sera célébrée le lundi, 25 juin, à 10 heures, en l'église St. Joseph, Esch- Alzette. Merci spécial pour les belles fleurs et les saintes messes. 4521</p>

Current quality in our collection



The workflow

1. Take a block of text
2. Identify the language
3. Calculate Confidence using dictionary
4. Run Tesseract OCR
5. Calculate new confidence
6. Replace existing ALTO block if new OCR is better

OCR software: Abbyy Finereader 8.1 in 2009

Dresden, 28. März. (Europapreß.) 3m Anschluß an
lkne Erweiböwsenversammlung, bie von anderen Erwerbs»
lösen gesprengt würbe, ereigneten sich gestern schwere Tu»
«nulte. Die Demonstranten zogen vor das Rathaus unb
vor bas Polizeipiäsibium, wo sie lmrch bie Polizei au«»
cinanbergeicW würben. Darauf versuchlen sie bie Kau^f.
laden in den umliegenden Straßen zu plunbern, jeboch
hatten bie «Inhaber bie Laben vorher geschlossen.

Statistics:

CLD2 language: GERMAN

60 words in total, 37 words in dictionary

347 characters in total, 196 chars in dictionary words

Confidence value : **56,48**

OCR software: Tesseract 3.04 in 2017

Dresden, 28. März. (Cuwpapreß.) Im Anschluß an
kne Erwerbglöfensockfammlng, die von anderen Erwerbs-
loscn gelptengt wurde, ereigneten sich gestern schwere Tu-
multe. Die Demonitrancken zogen vor das Rathaus squ
vor das Polizeipthium wo sie durch die Polizei aus-
einandergejcsgt wurden. Darauf versuchten sie die Kauf-
läden in den umliegenden Straßen zu plündern, jedoch
hatten die Inhaber die Läden vorher geschlossen.

Statistics:

CLD2 language: GERMAN

56 words in total, 44 words in dictionary

350 characters in total, 228 chars in dictionary words

Confidence value : **65,14 > 56,48**

HOCR output from Tesseract

```
<div class='ocr_page' id='page_1' title='image "P2_TB00012.tif"; bbox 0 0 1006 354;'>
  <div class='ocr_carea' id='block_1_1' title="bbox 0 0 1006 354">
    <p class='ocr_par' dir='ltr' id='par_1_1' title="bbox 0 0 1006 354">
      <span class='ocr_line' id='line_1_1' title="bbox 63 0 1006 46; baseline -0.006 -5">
        <span class='ocrx_word' id='word_1_1' title='bbox 63 6 218 46; x_wconf 75'
          lang='deu_frak' dir='ltr'>Dresden,</span>
        <span class='ocrx_word' id='word_1_2' title='bbox 241 5 290 40; x_wconf 82'
          lang='deu_frak' dir='ltr'>28.</span>
        <span class='ocrx_word' id='word_1_3' title='bbox 314 3 420 46; x_wconf 62'
          lang='deu_frak' dir='ltr'>Märs.</span>
        <span class='ocrx_word' id='word_1_4' title='bbox 445 1 686 45; x_wconf 74'
          lang='deu_frak' dir='ltr'>(Cuwpapreß.)</span>
        <span class='ocrx_word' id='word_1_5' title='bbox 713 1 768 38; x_wconf 75'
          lang='deu_frak' dir='ltr'>Jm</span>
        <span class='ocrx_word' id='word_1_6' title='bbox 791 0 942 43; x_wconf 73'
          lang='deu_frak' dir='ltr'>Ansichliuß</span>
        <span class='ocrx_word' id='word_1_7' title='bbox 966 9 1006 36; x_wconf 83'
          lang='deu_frak' dir='ltr'>an</span>
      </span>
    </p>
  </div>
</div>
```

Then:

1. Convert to ALTO, replace old TextBlock
2. Update METS file with new checksum

Result: choose best of 2

	Original OCR	New OCR	Best of both
#words	642628	631372	639109
#correctWords	455206	468608	512178
#characters	3567503	3604721	3572908
#correctCharacters	2316706	2324814	2635326
dictionary metric	64,94	64,49	73,76

Results for running OCR on 35 METS files

Challenges

- Dictionary metric can be improved upon
- OCR programs take a long time
- If the quality of scanned images is poor, OCR results are still bad in 2017

Thank you

Yves Maurer

yves.maurer@bnl.etat.lu

Bibliothèque nationale de Luxembourg

<https://github.com/ymaurer>