

Look, i can read now !

Koen Decorte



CD-Invest - Introduction

- International IBM i ISV and IBM business partner located in Antwerp, Belgium and Madrid, Spain
- Working with IBM i and its predecessors for more than 40 years
- Applications : CDQuery, CDSecure, CDView, CDLightning, CDReport, CDAccount and CDERP
- Expertise in RPG, SQL, Python, PHP, HTML, Unity, nodejs, linux, ...
- Website : www.cdinvest.eu
- Member of CEAC since 2018
- 6 IBM Champions in the company
- What others talk about, we do.

CD-Invest - Some of our customers



CD-Invest - IBM i Client Stories

Deknudt Frames

Building the framework for a thriving e-commerce operation with IBM i



ID-Logistics

Meeting the Challenges of a Pandemic with IBM i in the Cloud



JORI

Increasing Manufacturing Efficiency During COVID-19 With IBM i and advanced 3D-configurator



Diners Club Spain

Streamlining Customer Support with a Hybrid Cloud Application and IBM i



Wijnen Van Maele

Tracking wine production with blockchain on IBM i



Optimco

Introducing AI and a new customer experience in the car insurance industry on IBM i



CD-Invest - IBM i Client Stories

Fibrocity

Providing a comfortable seat with IBM i



Cras Woodgroup

Modernizing the wood industry with IBM i



Oris

Making vacations easier with IBM i



Steffimmo

Moving to IBM i on POWER9 in the cloud for growth



Stonetales properties

Upgrading and Centralizing on the Cloud with IBM i



Winsol

Digitizing manufacturing on IBM i



CD-Invest - IBM i Client Stories

CSM

Empower more small
businesses to access
global trade



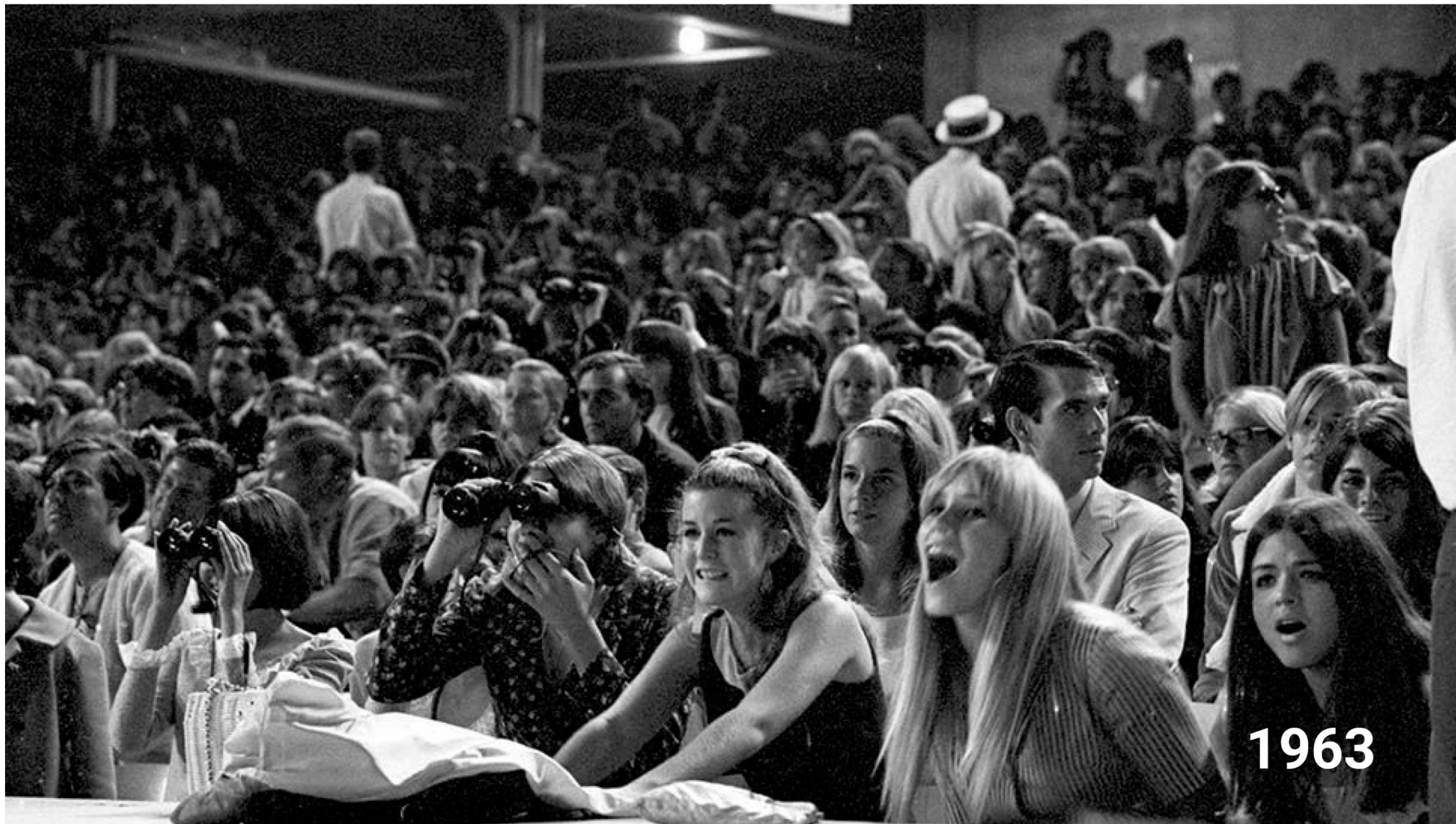
Bonehill

Adapting IBM i to the
modern web



Read more on on <https://www.ibm.com/it-infrastructure/us-en/resources/power/ibm-i-customer-stories/>

Introduction



1963



2015

Importance of pictures

- 3.2 billion pictures are shared online each day
- They contain brands, products, services, ads ...
Interesting to find out if they are yours.
- 80% of those images don't mention brand or product in the accompanying text.

→ Insight missing and what about reading or understanding the text ??

But first things first

What is AI ?

Artificial Intelligence (AI)

= science of making computers do things that require intelligence when done by humans.

Above intelligent or Abysmal idiot ?

What is OCR?

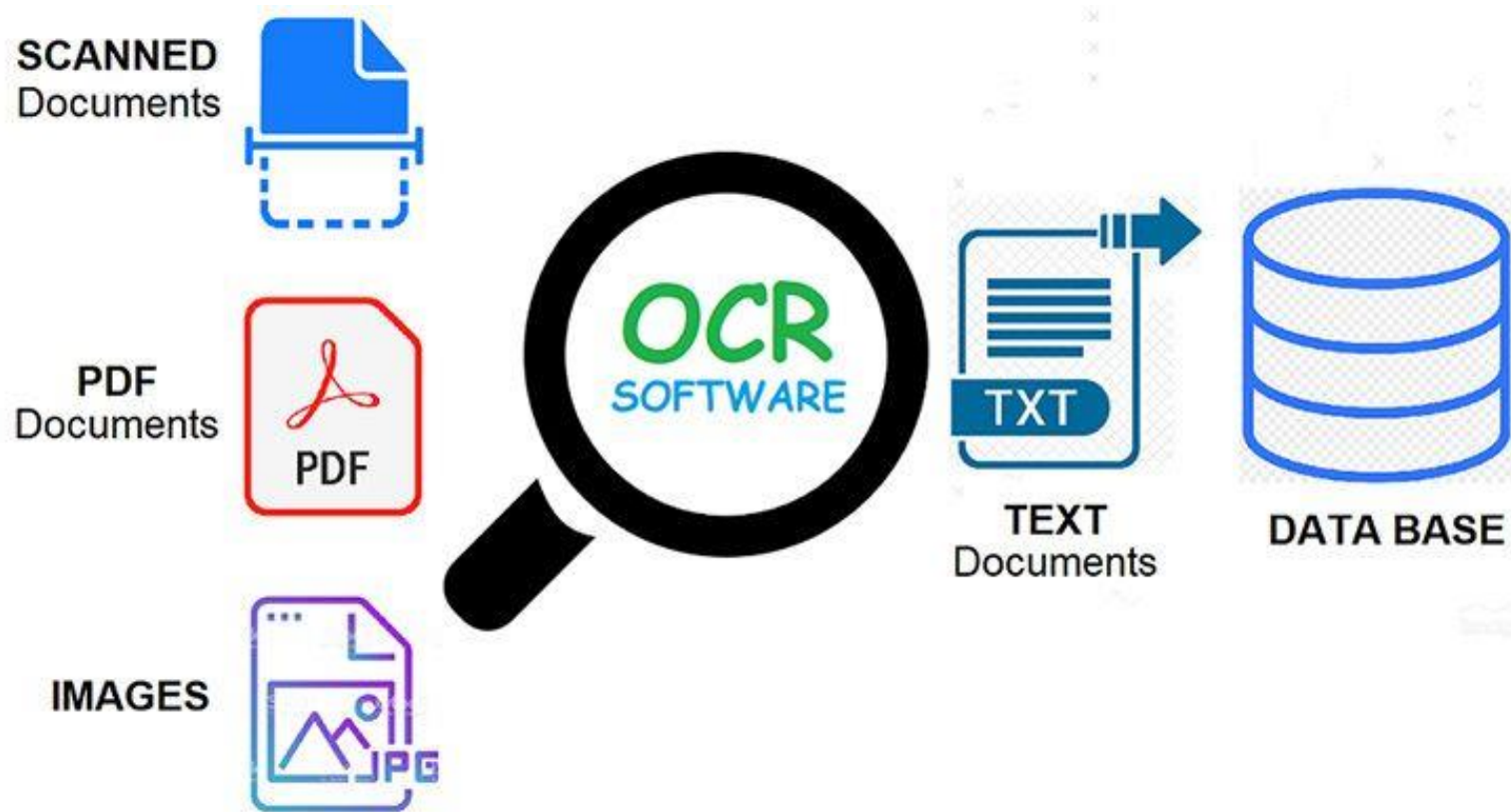
- Stands for Optical Character Recognition
- Extracts the text from a given image

My invention relates to statistical machines
of the type in which successive comparisons
are made between a character and a charac-



My invention relates to statistical machines
of the type in which successive comparisons
are made between a character and a charac-

What is OCR ?



What is OCR ?

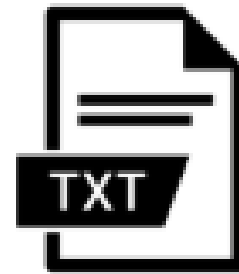
OCR predates electronic computers !!



Scanned Document
(Image)



OCR

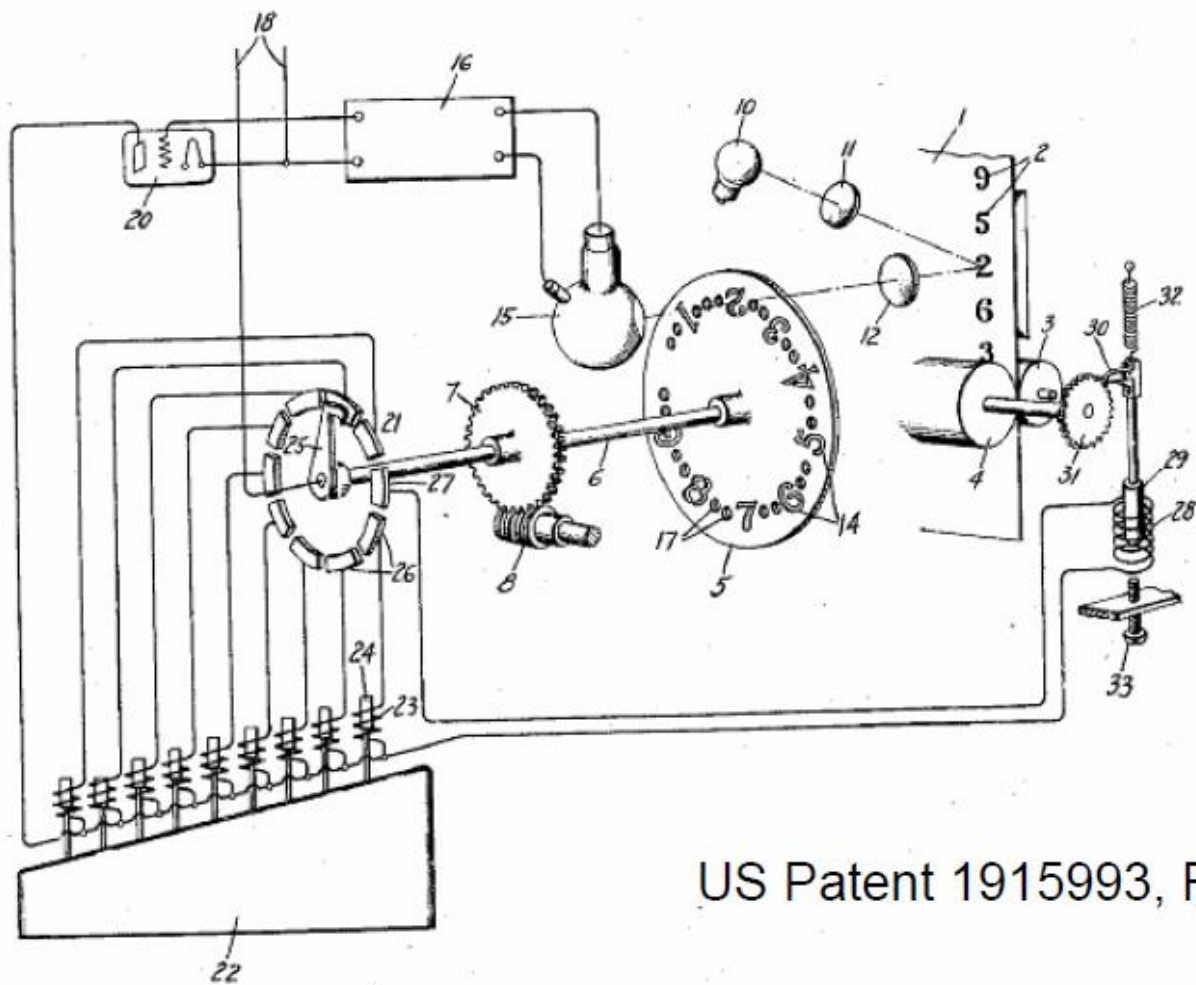


Text

What is OCR?

- Invented by Gustav Tauschek
- Tauschek obtained a patent on OCR
 - 1929 in Germany
 - 1935 in USA
- Tauschek's machine
 - Was a mechanical device
 - Uses templates, light and photodetector
 - When a light was directed towards the templates no light reach the photodetector

What is OCR?



US Patent 1915993, Filed Apr 27, 1931

OCR history

- 1929 – Digit recognition machine
- 1953 – Alphanumeric recognition machine
- 1965 – US mail sorting
- 1965 – British banking system
- 1976 – Kurzweil reading machine
- 1985 – Hardware assisted PC software
- 1988 – Software-only PC software
- 1994-2000 – Industry consolidation

What is OCR ?

- OCR Subprocesses
 - Preprocessing of the Image
 - Text Localization
 - Character Segmentation
 - Character Recognition
 - Post Processing

OCR use cases

- building license plate readers
- digitizing invoices
- reading container numbers
- digitizing ID cards
- digitizing letters, insurance documents, ...

Project Tesseract

Project Tesseract

- History of Tesseract
 - Open source OCR engine
 - Developed by HP between 1985 and 1995
 - Never used in an HP product
 - Rated highly at The Fourth Annual Test of OCR Accuracy in 1995
 - In 2005 HP transferred Tesseract to the ISRI and released it as open source
 - ISRI == Information Science Research Institute
 - The development is currently led by Google (since 2006 !)

Project Tesseract

- Tesseract is an OCR Engine and is NOT a complete OCR program
 - Originally intended to serve as a component part of other programs
 - Works from the command line
 - Has no GUI
 - Integration to many programming languages
 - Runs on IBM i



● Tesseract OCR
Search term

● OCRopus
Search term

● Ocular OCR
Search term

● SwiftOCR
Search term



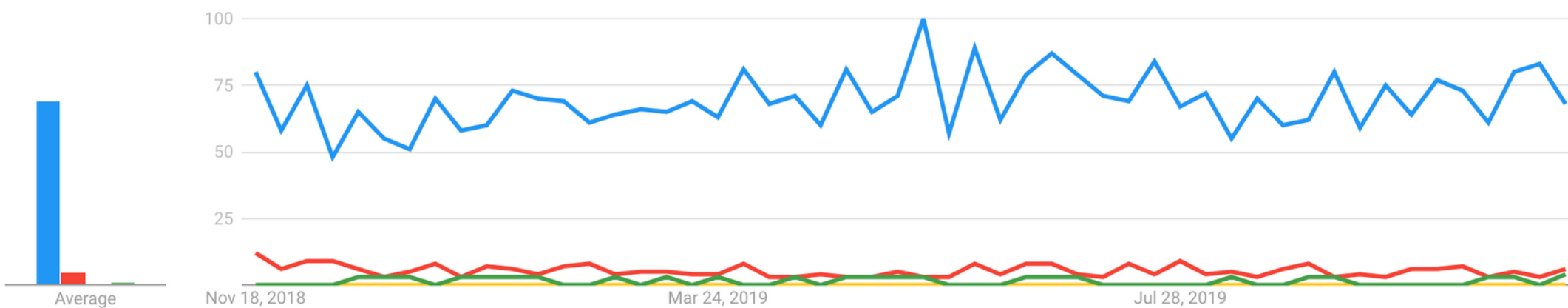
Worldwide ▼

Past 12 months ▼

All categories ▼

Web Search ▼

Interest over time 



Install Tesseract

Install tesseract

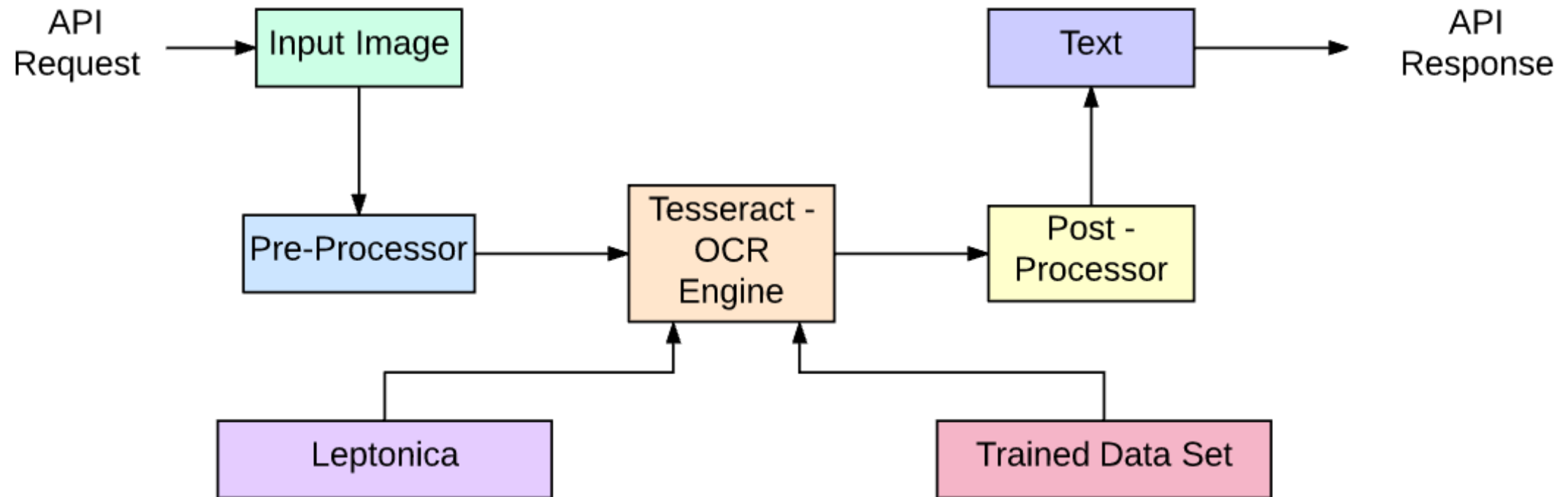
- Install from open source package manager (if needed I can provide compile scripts since I've given mine to IBM).

leptonica-devel	1.80.0-2	@ibmi-base
leptonica-tools	1.80.0-2	@ibmi-base

tesseract-devel	4.1.1-1	@ibmi-base
tesseract-tessdata	4.1.1-1	@ibmi-base
tesseract-tools	4.1.1-1	@ibmi-base

How does Tesseract work ?

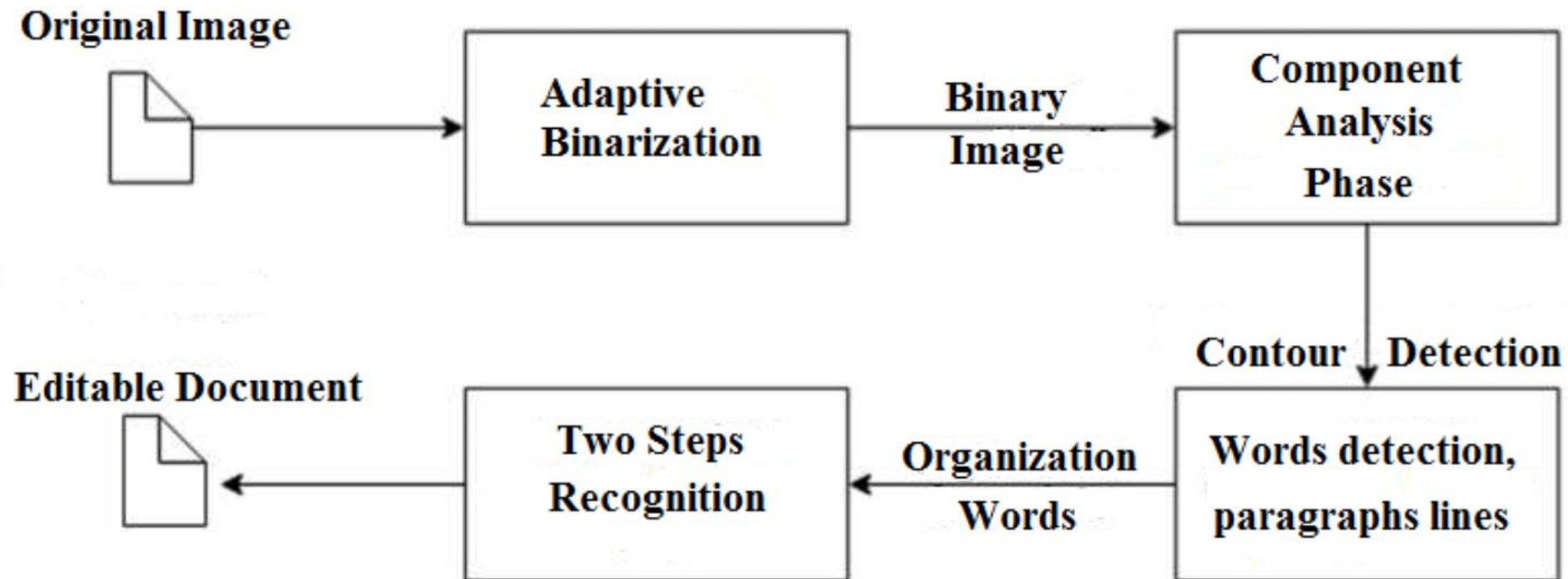
OCR Process Flow



How does Tesseract work ?

- Tesseract 4.00 includes a new neural network subsystem configured as a text line recognizer.
- It has its origins in OCRopus Python LSTM implementation but has been redesigned for **Tesseract in C++**. **The neural network system in Tesseract pre-dates TensorFlow.**
- To recognize an image containing a single character, we typically use a Convolutional Neural Network (CNN). Text of arbitrary length is a sequence of characters, and such problems are solved using RNNs and LSTM (Long short-term memory) is a popular form of RNN (recurrent neural network).

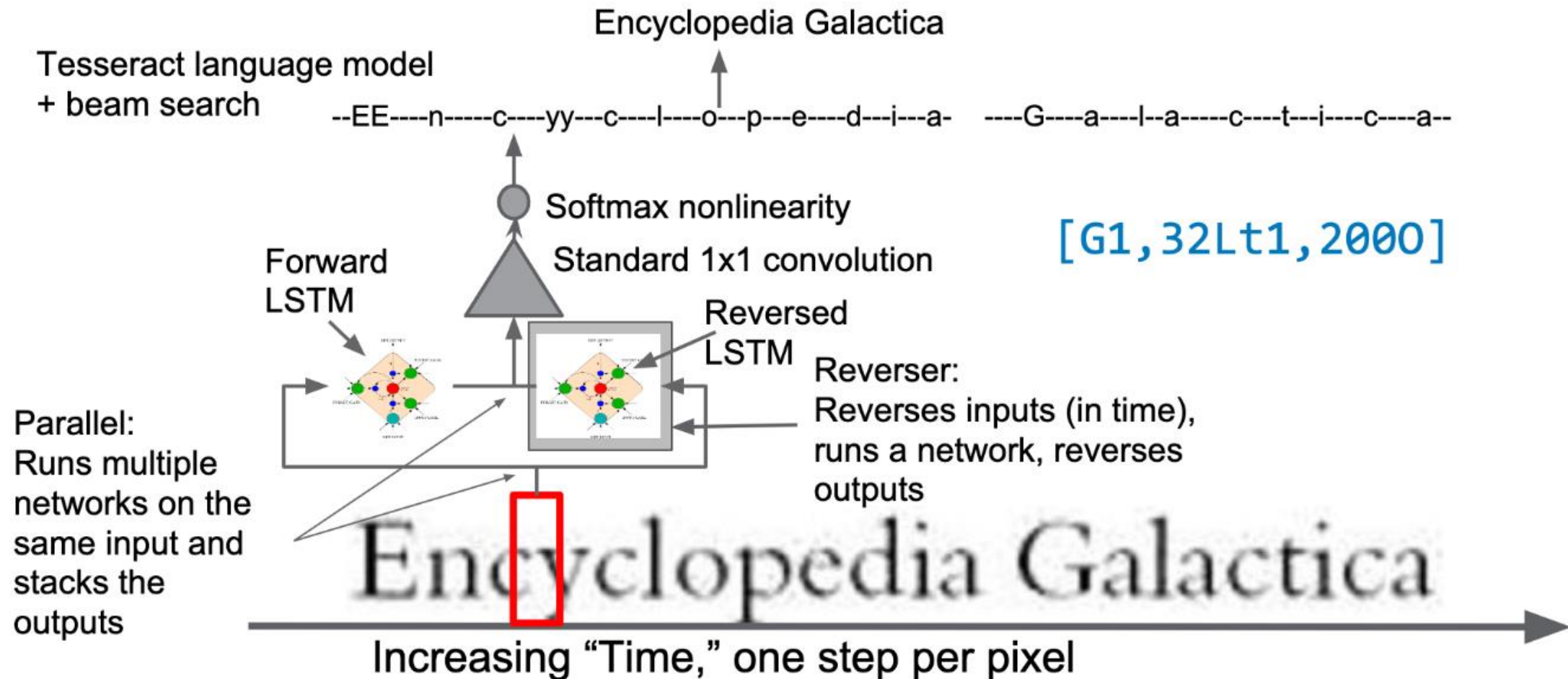
How does tesseract work ?



How Tesseract Works? (Old way and many traditional OCR

1. Adaptive thresholding on the input image
2. Analyze connected components in the binary image
3. Find text lines and words
4. First pass of recognition process
 - Attempts to recognize each word in turn
5. Satisfactory words are passed to adaptive trainer
6. Lessons learned are employed in a second pass
 - Used for words not satisfactory recognized
7. Producing the output text

How Tesseract uses LSTMs...



Running Tesseract ?

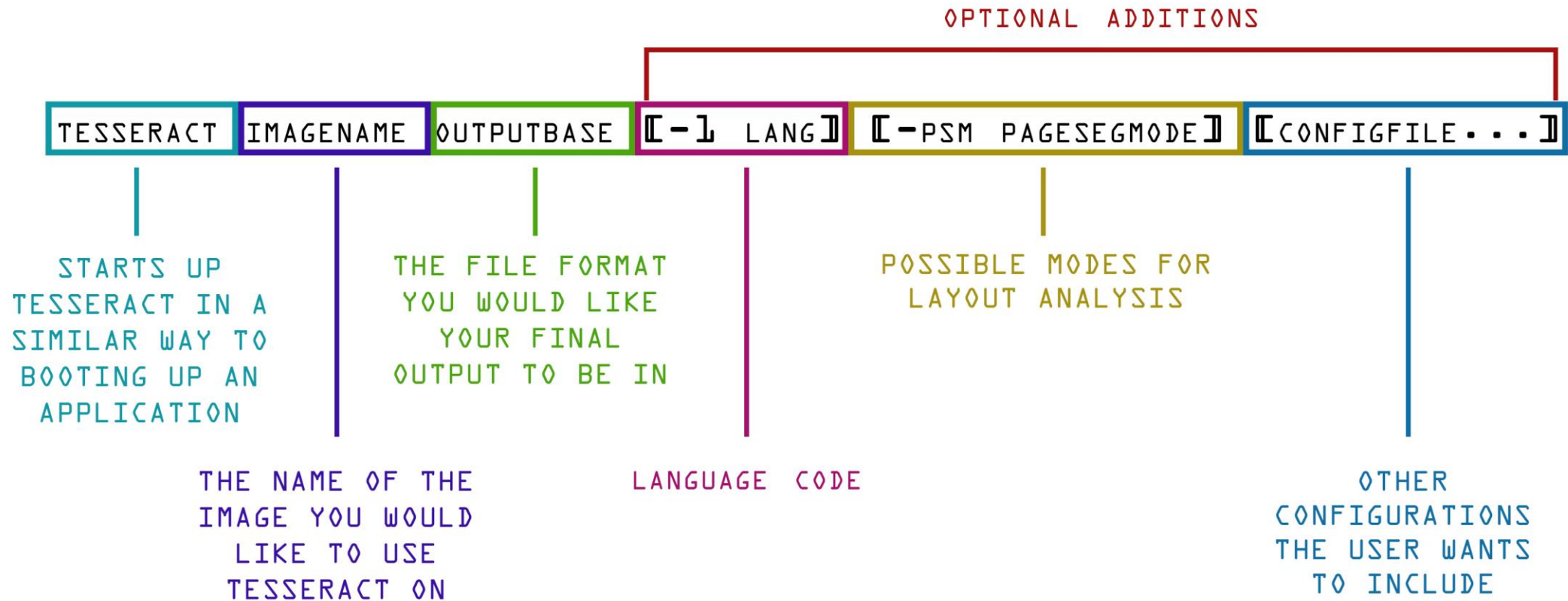
Running tesseract

```
$  
> tesseract  
Usage:  
  /QOpenSys/pkgs/bin/tesseract --help | --help-extra | --version  
  /QOpenSys/pkgs/bin/tesseract --list-langs  
  /QOpenSys/pkgs/bin/tesseract imagename outputbase [options...] [configfile...]  
  
OCR options:  
  -l LANG[+LANG]          Specify language(s) used for OCR.  
NOTE: These options must occur before any configfile.  
  
Single options:  
  --help                  Show this help message.  
  --help-extra            Show extra help for advanced users.  
  --version               Show version information.  
  --list-langs            List available languages for tesseract engine.  
$
```

Running tesseract

```
> tesseract --version  
tesseract 4.1.1  
  leptonica-1.80.0  
    libgif 5.1.4 : libpng 1.6.37 : libtiff 4.0.9 : zlib 1.2.13 : libwebp 1.0.2  
$
```

Running Tesseract



Running Tesseract

```
$ tesseract image_path stdout
```

To write the output text in a file:

```
$ tesseract image_path text_result.txt
```

To specify the language model name, write language shortcut after -l flag, by default it takes English language:

```
$ tesseract image_path text_result.txt -l eng
```

Running tesseract

receipt-home-depot.jpg

receipt-mels.jpg

receipt-sunshine.jpg

receipt-tadish.jpg

```

Welcome to Mel's

Check #: 0001                      12/20/11
Server: Jcsh F                     4:38 PM
Table: 7/1                         Guests: 2
-----
2 Beef Burgr (@9.95/ea)            19.90
  SIDE: Fries
1 Bud Light                        3.79
1 Bud                              4.50
-----
Sub-total                          28.19
Sales Tax                          2.50
TOTAL                             30.69
-----
Balance Due                        30.69

Thank you for your patronage!
```

1425788830...depot.jpg.txt
1425788831...t-mels.jpg.txt
1425788834...shine.jpg.txt
1425788834...tadish.jpg.txt

welcome to MeT's

```

Check #: 0001 12/20/11
Server: Josh F 4:38 PM
Table: 7/1 Guests: 2
2 Beef Burgr (@9.95/ea) 19.90
SIDE: Fries
1 Bud Light 3.79
1 Bud 4.50
Sub-total 28.19
SaTes Tax _____mg.$g
TOTAL 30.69
eai;r};;"13;; "
"Tédfee

Thank you for your patronage!
```

Running Tesseract

```
> tesseract --list-langs  
List of available languages (6):  
bel  
deu  
eng  
fra  
ita  
spa  
$
```

Running Tesseract

```
$ tesseract image_path text_result.txt -l eng --psm 6
```

There is also one more important argument, OCR engine mode (oem). Tesseract 4 has two OCR engines — Legacy Tesseract engine and LSTM engine. There are four modes of operation chosen using the --oem option.

- 0 Legacy engine only.
- 1 Neural nets LSTM engine only.
- 2 Legacy + LSTM engines.
- 3 Default, based on what is available.

Running Tesseract

- Beware of memory / CPU limitations using LSTM ! It is CPU intensive.
- In addition, LSTM networks have pre-trained models that can be used in OCR enhancement. You can also adapt these models to the specific requirements of OCR with little effort.
- In practice, LSTM networks are mainly used to recognize handwritten texts and old documents with high accuracy - where OCR algorithms have problems.

Running Tesseract – supported files

- JPG
- PNG
- GIF
- PNM
- TIFF
- Unfortunately, Tesseract engine can't read PDF file. (Use ghostscript or imagemagick to convert PDF to images).

Running Tesseract - Programming Language Package

- Gosseract (Golang)
- Tess4J (Java)
- RTesseract (Ruby)
- PyTesseract (Python)
- Tesseract.js (Javascript)
- Thiagoalessio TesseractOCR (PHP)
- etc...

Running Tesseract - Tesseract hOCR output format

- either text (txt) or hOCR, an html-format with embedded segmentation info
- bounding box: x0y0 x1y1
- `<span class='ocrx_word' id='word_1_33' title='bbox 1584 1199 1997 1284; \`
- `x_wconf 87' lang='deu-frak' dir='ltr'>Verhältnisse.`
- hOCR has word tokens (separated by white space) as smallest unit
- Ben Kiessling (Nidaba project), Kay Würzner have achieved character xml output

Running Tesseract

```
<meta name='ocr-system' content='tesseract 4.1.1' />
<meta name='ocr-capabilities' content='ocr_page ocr_carea ocr_par ocr_line ocrx_word ocrp_wconf' />
</head>
<body>
<div class='ocr_page' id='page_1' title='image "invoice-1.jpg"; bbox 0 0 595 842; ppageno 0'>
<div class='ocr_carea' id='block_1_1' title="bbox 26 25 145 56">
<p class='ocr_par' id='par_1_1' lang='fra' title="bbox 26 25 145 56">
<span class='ocr_line' id='line_1_1' title="bbox 27 25 145 40; baseline 0.008 -1; x_size 19.26087; x_descenders 5.2608695; x_ascenders 3">
<span class='ocrx_word' id='word_1_1' title='bbox 27 25 145 40; x_wconf 92'>Prestatiestaat</span>
</span>
<span class='ocr_line' id='line_1_2' title="bbox 26 47 63 56; baseline 0 0; x_size 22.5; x_descenders 5.5; x_ascenders 5.5">
<span class='ocrx_word' id='word_1_2' title='bbox 26 47 63 56; x_wconf 54'>core</span>
</span>
</p>
</div>
```

Running Tesseract – tsv output

level	page_num	block_num	par_num	line_num	word_num	left	top	width	height	conf	text
1	1	0	0	0	0	595	842	-1			
2	1	1	0	0	0	26	25	119	31	-1	
3	1	1	1	0	0	26	25	119	31	-1	
4	1	1	1	1	0	27	25	118	15	-1	
5	1	1	1	1	1	27	25	118	15	92	<u>Prestatiestaat</u>
4	1	1	1	2	0	26	47	37	9	-1	
5	1	1	1	2	1	26	47	37	9	54	core
2	1	2	0	0	0	26	71	76	11	-1	
3	1	2	1	0	0	26	71	76	11	-1	
4	1	2	1	1	0	26	71	76	11	-1	
5	1	2	1	1	1	29	61	33	28	47	<u>ringe</u>
5	1	2	1	1	2	66	61	36	28	83	0172001
2	1	3	0	0	0	22	88	551	5	-1	
3	1	3	1	0	0	22	88	551	5	-1	
4	1	3	1	1	0	22	88	551	5	-1	
5	1	3	1	1	1	22	88	551	5	95	
2	1	4	0	0	0	31	103	125	8	-1	
3	1	4	1	0	0	31	103	125	8	-1	
4	1	4	1	1	0	31	103	125	8	-1	
5	1	4	1	1	1	31	103	52	8	25	<u>Nedocunet:</u>
5	1	4	1	1	2	99	103	57	8	42	2021010030
2	1	5	0	0	0	22	120	551	5	-1	
3	1	5	1	0	0	22	120	551	5	-1	
4	1	5	1	1	0	22	120	551	5	-1	
5	1	5	1	1	1	22	120	551	5	95	
2	1	6	0	0	0	22	133	551	5	-1	
3	1	6	1	0	0	22	133	551	5	-1	

Tesseract Training

- enable recognition of a new "language"
- rather, what is trained are glyph shapes:
 - new alphabet (Latin, Cyrillic, Greek etc.)
 - new typeface (Antiqua, Fraktur)
 - new font (special instance of a typeface, e.g. 12 pt Caslon italic)
- optionally add language data (wordlists)

Tesseract Training

- training data are in files of the form LANG.traineddata
- glyph shape training data and language support data (wordlists) are tied up in the same file
- language data can be exchanged without retraining (better and larger wordlists)

Tesseract accuracy

268702 Characters

5954 Errors

97.78% Accuracy

Errors	Marked	Correct-Generated
--------	--------	-------------------

372	0	{ }-{ }
-----	---	---------

246	0	{ü }-{ii }
-----	---	------------

228	0	{ }-{ }
-----	---	---------

215	0	{I }-{J }
-----	---	-----------

211	0	{v }-{V }
-----	---	-----------

208	0	{ }-{<\n> }
-----	---	-------------

175	0	{u }-{n }
-----	---	-----------

140	0	{n }-{u }
-----	---	-----------

136	0	{c }-{e }
-----	---	-----------

Tesseract guidelines

- Tesseract performs well when document images follow the next guidelines:
- Clean segmentation of the foreground text from background
- Horizontally aligned and scaled appropriately
- High-quality image without blurriness and noise

Tesseract guidelines

- The latest release of Tesseract 4.0 supports deep learning based OCR that is significantly more accurate. The OCR engine itself is built on a Long Short-Term Memory (LSTM) network, a kind of Recurrent Neural Network (RNN).
- Tesseract is perfect for scanning clean documents and comes with pretty high accuracy and font variability since its training was comprehensive. I would say that Tesseract is a go-to tool if your task is scanning of books, documents and printed text on a clean white background.

Application examples using
Tesseract ?

Tesseract examples

- **Healthcare industry**

As a rule, hospitals and doctors' offices keep medical records in written form. In large quantities, these are therefore difficult to search. Tesseract can digitize these records, organize them - and thus make them easily searchable. Doctors and nurses can thus automatically analyze large volumes of medical records and extract important information. This leads to more efficient diagnosis and treatment of patients.

- **Finance**

Financial documents such as bank statements, Invoices and tax returns are still often created in writing. Searching these is therefore time-consuming. Tesseract can index and categorize these documents quickly and automatically. Banks can thus automatically read in checks, for example, and thus significantly reduce the manual workload.

Tesseract examples

- **Logistics**

In the logistics industry, it is important to be able to quickly access information such as package numbers, inventory figures and shipping addresses. Tesseract enables automatic recognition of product labels and Barcodes. This leads to faster and more accurate recording of inventories. In this way, companies in logistics can increase their efficiency and avoid bottlenecks in inventory management.

- **Mobile applications**

Tesseract can be embedded as a component in mobile apps to recognize text within images on mobile devices. This is particularly useful for applications such as translation and text recognition apps.



Browse : /invoice2data/invoice2data.sh

Record : 18 of 46 by 18

Column : 1 204 by 131

Control :

```
.....1.....2.....3.....4.....5.....6.....7.....8.....9.....0.....1.....2.....3.
#Execute invoice2data inside chroot container

chroot /Q0penSys/invoice2data invoice2data --input-reader pdftotext --template-folder ocr/templates ocr/pdfs/$invoice --output-f

#Copy the result back to ifs

cp /Q0penSys/invoice2data/ocr/json/$json /invoice2data/json/$json 2>>/invoice2data/log.txt

#move the invoice to the done folder

mv /invoice2data/invoices/$invoice /invoice2data/done/$invoice 2>>/invoice2data/log.txt

#delete files inside de chroot container

rm /Q0penSys/invoice2data/ocr/pdfs/$invoice 2>>/invoice2data/log.txt
rm /Q0penSys/invoice2data/ocr/json/$json 2>>/invoice2data/log.txt

echo "Done Processing -----" >>log.txt
}
```

F3=Exit F10=Display Hex F12=Cancel F15=Services F16=Repeat find F19=Left F20=Right

1 records folded.

MA

A

MW

03/012



Browse : /invoice2data/json/coenen.json

Record : 1 of 13 by 18

Column : 1 59 by 131

Control :

.....1.....2.....3.....4.....5.....6.....7.....8.....9.....0.....1.....2.....3.

*****Beginning of data*****

```
[
  {
    "account": 70202,
    "amount": 2736.35,
    "baseamount": "2.299.45",
    "currency": "EUR",
    "date": "2022-02-24",
    "desc": "Invoice from Coenen Neuss GmbH & Co. KG",
    "invoice_number": "40201083/ 1",
    "issuer": "Coenen Neuss GmbH & Co. KG",
    "vat": "436.90"
  }
]
```

*****End of Data*****

F3=Exit F10=Display Hex F12=Cancel F15=Services F16=Repeat find F19=Left F20=Right

1 records folded.

MA

A

MW

03/012



Browse : /invoice2data/json/Kylma_2.json

Record : 1 of 30 by 18

Column : 1 59 by 131

Control :

.....1.....2.....3.....4.....5.....6.....7.....8.....9.....0.....1.....2.....3.....

*****Beginning of data*****

```
[
  {
    "amount": 1032.62,
    "currency": "EUR",
    "date": "2023-03-14",
    "desc": "Invoice from Kylma AB",
    "invoice_number": "3213036-1",
    "issuer": "Kylma AB",
    "line": [
      [
        "916868",
        "2023-03-28",
        "1,00",
        "52,54"
      ],
      [
        "590240",
        "2023-03-28",
```

F3=Exit F10=Display Hex F12=Cancel F15=Services F16=Repeat find F19=Left F20=Right

1 records folded.



MA

A


MW


03/012


Ingave aankoopfacturen

 Einde ingave  Bevestigen


 PDF

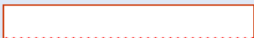
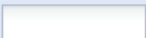
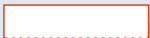
 Convert Pdf to Js

Datum ingave: 11.09.2023 


Journaal: AF - AANKOOPFACTUREN 

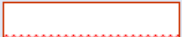

DocumentNr: 2300202

Leverancier: 

Factuurnummer:  Datum:  Vervaldatum: 


Blokken: ☒ Betaald: ☐

Valutadatum: 11.09.2023 

Te betalen bedrag:  EUR Btw: 4 - 21% Soort:  Land: BE

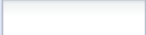
Totaal zonder btw: 0

Btw aftrekbaar: 0

Btw niet aftrekbaar: 0 Rekening: 

Te betalen btw: 0


Investering: 0

Disconto: 0 Vervaldatum: 


Buiten maatstaf: 0

Kredietbeperking: 0

Te betalen leeggoed: 0

 Debet: EUR 0,00 Credit: EUR 0,00

Te betalen bedrag niet gelijk aan de som van de uitsplitsing


	K/L/R	Rekening	Rekening omschrijving	Kostenplaats	Kostenplaats omschrijving	Bedrag EUR	D/C	Btw vak	Omschrijving	Personeel	Personeel omschrijving
	L	440000				0,00	C				
	R					0,00	D				


Ingave aankoopfacturen

 Einde ingave

 Bevestigen

 PDF

Datum ingave: 11.09.2023 

Journaal: AF - AANKOOPFACTUREN 

DocumentNr: 2300202

Leverancier: /accounting/scan/

Factuurnummer:

Blokkeren:

Te betalen bedrag:

Totaal zonder btw:

Btw aftrekbaar:

Btw niet aftrekbaar:

Te betalen btw:

Investering:


Disconto:

Buiten maatstaf:

Kredietbeperking:


Te betalen leeggoed:

	Filename	Creation date	Creation time
<input type="checkbox"/>	20230824153541.pdf	24.08.2023	15.33.03
<input type="checkbox"/>	20230906130652.pdf	06.09.2023	13.23.18
<input checked="" type="checkbox"/>	820248_2023-08-15_FR_E-INVOICE_23-631304081-002.pdf	24.08.2023	15.35.37
<input type="checkbox"/>	ZZZ_0000001.pdf	24.08.2023	15.24.40
<input type="checkbox"/>	ZZZ_0000002.pdf	24.08.2023	15.26.26

 Confirm

 Cancel

Te betalen bedrag niet gelijk aan de s

	K/L/R	Rekening	Rekening omschrijving	Kostenplaats	Kostenplaats omschrijving	Bedrag EUR	D/C	Btw vak	Omschrijving	Personeel
	L	440000				0,00	C			
	R					0,00	D			

Ingave aankoopfacturen

Einde ingave Bevestigen

Datum ingave: 11.09.2023

Journaal: AF - AANKOOPFACTUREN DocumentNr: 2300202

Leverancier: 48348 - DKV NIEDERLASSUNG DUESSELDORF - DUESSELDORF

Factuurnummer: 1234567 Datum: 11.09.2023 Vervaldatum: 11.09.2023

Blokkeren: ☒ Betaald: ☐ Valutadatum: 11.09.2023

Te betalen bedrag: 2100 EUR Btw: 4 - 20% Land: FR

Totaal zonder btw: 2100

Btw niet aftrekbaar: 420 Rekening: 411020 - TERUG TE VORDEREN BTW FRANKRIJK

Disconto: 0 Vervaldatum: 11.09.2023

Buiten maatstaf: 0

Kredietbeperking: 0

Te betalen leeggoed: 0

PDF

Convert Pdf to Jsou

1 of 1

Automatic Zoom

E-FACTURE /
E-REKENING

Pour livraisons et prestations en
Voor dienstverlening en leveringen in
France

Adresse du client/adresse de facturation
D.I.C. INVIA
Deknudt Invest en Consult
Charline Bourgois
Breestraat 31A
8540 DIEERLIJK
BELGIEN

Número client /
Klantnummer:
0000820248

numéro fiscal EU client/ No. fiscal national /
Kl. btw Id./Nat. bel. nr.:
BE0894259628

Número de factura /
Factuurnummer:
234631304081002

Date de facturation /
Factuurdatum:
15.08.2023

Page / Pagina: 1 / 1

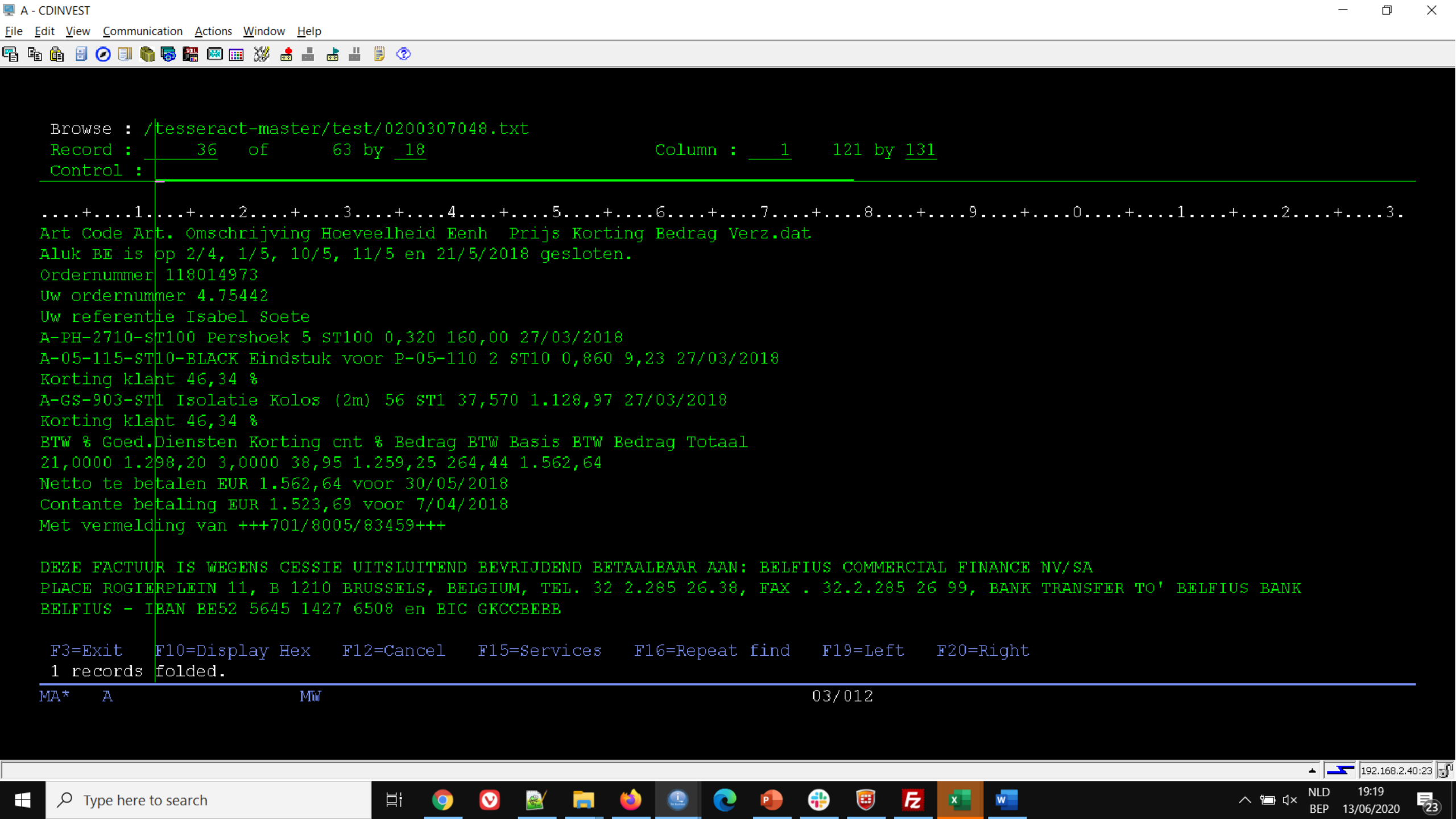
Date de livraison	Station service	Número de la transaction	Service	Kilo-métrage	Produit	Unité	Quantité	Prix unitaire TTC	Produit net	Produit TTC	Scalé de livraison	Scalé de service	Produit net	Produit TTC	Scalé de livraison	Scalé de service	Produit net	Produit TTC
10.08.2023	ESGO	1301647	ESGO	10.47	2 (LITRE 95 (SUPER))	0036	118	24.980	2.0480	1.8035	42.57	1.17	43.74	8.75	52.49			
TOTAL: 24.980 42.57 1.17 43.74 8.75 52.49																		
Résumé des sommes totales par type de prestation avec 20.00% TVA																		
EURO 95 (SUPER) 0036 118 24.980 42.57 1.17 43.74 8.75 52.49																		
EURO 95 (SUPER) 0036 118 24.980 42.57 1.17 43.74 8.75 52.49																		

Statistiques TVA

20.00 %	43.74	8.75	52.49
	43.74	8.75	52.49

Nombre de TVA DKV : 1875427548025

Debet: EUR 840,00 Credit: EUR 2.100,00					Totaal debet niet gelijk aan totaal credit					
K/L/R	Rekening	Rekening omschrijving	Kostenplaats	Kostenplaats omschrijving	Bedrag EUR	D/C	Btw vak	Omschrijving	Personeel	Personeel omschrijving
L	440000	Leveranciers			2.100,00	C				
R	411020	Terug te vorderen btw Frankrijk			420,00	D		DKV NIEDERLASSUNG DUESSELDORF		
R	411003	Terug te vorderen btw België medecontractant			420,00	D		DKV NIEDERLASSUNG DUESSELDORF		





Batch Content

Import_Winsol_Invoice

1: Winsol Invoice
Page 1
Page 2

2: Winsol Invoice
Page 1
Page 2

3: Winsol Invoice
Page 1
Page 2

4: Winsol Invoice
Page 1
Page 2

5: Winsol Invoice
Page 1
Page 2
Page 3
Page 4

6: Winsol Invoice
Page 1
Page 2
Page 3
Page 4

7: Winsol Invoice
Page 1
Page 2

8: Winsol Invoice
Page 1
Page 2
Page 3
Page 4

9: Winsol Invoice
Page 1
Page 2

10: Winsol Invoice
Page 1
Page 2

11: Winsol Invoice
Page 1
Page 2

12: Winsol Invoice
Page 1
Page 2

Classification Result

The document was classified to the following class: Invoice

40 fields valid, 1 field invalid (24 fields invisible, 2 fields read-only)

1 Pagina 1 (1 invalid field) 2 Pagina 2 (0 invalid fields)

Streepjescode groep [v. 5.0.4]

Streepjescode: 0200601110

Firma groep

Firma identificatie: AC

Firmanaam: AC - WINSOL ACTUELL N.V.

BTW-nummer: BE0448902241

Leveranciersgroep

Leveranciersidentificatie: 00002083

Leveranciersnaam: ALUK BELGIUM

BTW-nummer: BE0888805159

Rekeningnummer:

IBAN nummer: BE29001853967464

ACCOUNT: 600011

ANALYTICALACCOUNT1: 4000

Algemene groep

Factuurtype: I

Factuurnummer: 7020006494

Factuurdatum: 8/06/2020

InPlace Editor - Field "INVOICENUMBER"

7020006494

INVOICENUMBER

7020006494

Current Error

Please confirm this field content.

ALUK

ALUK Belgium NV
Zwaarveld 44
B-9220 Hamme

BTW-nr: BE0888805159
RPR: RPR. DENDERMONDE
Telefoon: +32 52 48 48 48
Fax: +32 52 48 48 16
E-mail: info.be@aluk.com

Factuur
7020006494

Klantennummer: 000056
Winsol Actuel NV
Afdeling Aluminium
Roosdaanstraat 542
B-8870 Izegem

BTW-nr: BE0448902241
Fax: +32 51 33 19 91

9 9 JUNI 2020

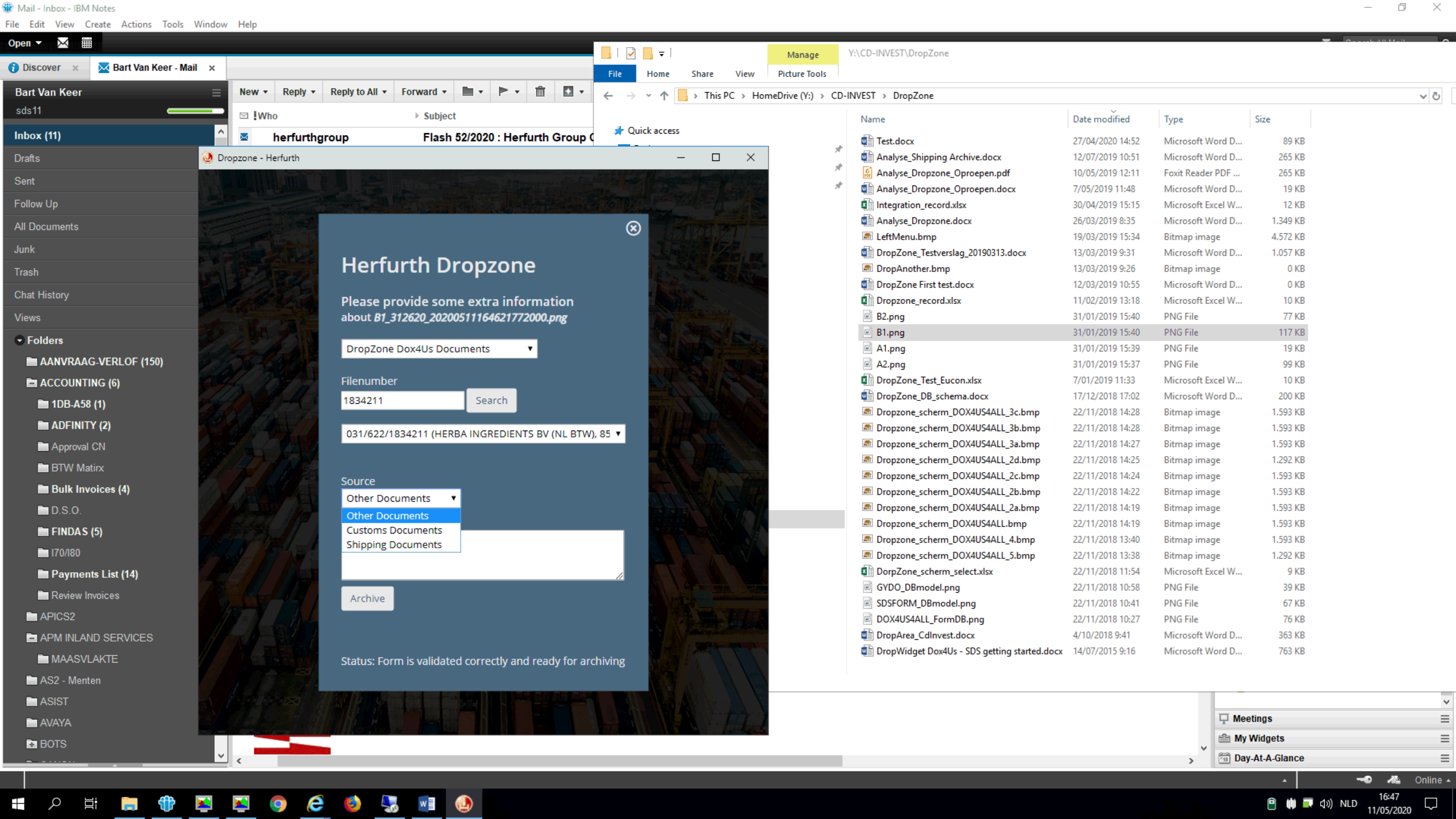
Factuurdatum: 8/06/2020 Vervaldag: 29/08/2020 Pag. 1/1

Art. Code	Art. Omschrijving	Hoeveelheid	Eenh.	Prijs	Korting	Bedrag	Verz.dat
Betsing van deze factuur is te voldoen op rekening: BNP PARIBAS FORIS - IBAN: BE29 0018 5396 7464 - BIC: CEBASE33							
Ordernummer: 120019302 Uw ordernummer: 12 04.0316.039.4.136173							
Uw rchivante: Isabel Soete							
P-03-P211-9005PA-6-500	Verslagtoefn vln 54cm	2	PCS	14,010		182.13	0/09/2020
P-05-039-9005PA-6-500	Gastat 39cm	2	PCS	4,399		57.97	0/09/2020

BTW %	Goed Diensten	Korting cfr. %	Bedrag	BTW Basis	BTW Bedrag	Totaal
21.0000	239.20	3.0000	7.16	232.02	48.72	287.92
Netto te betalen: EUR 207.92 voor 29/08/2020						
Contante betaling: EUR 200.74 voor 19/06/2020						
Met vermelding van: ***7020006494***						

BNP PARIBAS FORIS - IBAN: BE29 0018 5396 7464 - BIC: CEBASE33

ALUK Belgium NV - Zwaarveld 44 - 9220 Hamme
T: +32 52 48 48 48 F: +32 52 48 48 16





Mail2Me

Download PDF

Print PDF

**HERFURTH LOGISTICS**

Operations
 Station 1040
 22 rue de la Gare 1040
 +32 3 27 08 17

Administration
 ACCOUNTING
 22 rue de la Gare 1040
 +32 3 27 08 18

HERFURTH LOGISTICS NV
 Registered Office: CASSIERSTRAAT 19
 2000 ANTWERPEN 6, BELGIUM

VAT: BE0445364909
 R.P.R.: 0445.364.909, ANTWERPEN,
 AFDELING ANTWERPEN
 Registration nr. customs representative 2005

Invoice Nr. 027/908601 Date 09/12/2019

STEPAN DEUTSCHLAND GMBH
 RODENKIRCHENER STR. 400
 50389 WESSELING
 GERMANY

ORIGINAL

Please mention on payment: 027/908601
 Due date: 09/12/2019

Your customer nr.: S16359
 Your VAT nr.: 06811981536

Bank: IBAN: BIC/SWIFT
 KBC: BE3410064328197 FREDERIX EUR
 BNP PARIBAS FORIS: BE43220066330008 GERARDUS EUR

Invoice total amount			
Base VAT excl.	725,00 EUR		
VAT 0%	0,00 EUR	(on 725,00 EUR)(1)	
In our favour	725,00 EUR		

(1): Transport exemption from VAT: Art 146, a) from Eur. VAT Council Directive

Vessel/Voyage	Monte Alegre/ 9485	Our File	031/610/832411
Port of Loading	Antwerpen, BE	Total Packages	1
Port of Discharge	Ei Iskandariya (* Alexandria), EG	Total Weight	286,00 kg
E.T.S.	01/12/2019		
E.T.A.	14/12/2019		

Goods: harmless chemicals

Description	Quantity	Base	Amount	Curr.	Rate of Exchange	VAT%	Amount in EUR VAT Class
Compuem cool Fee to Fed Mailport	1,00	535,00	535,00+	EUR	1,000000	0%	535,00+ (1)
Our handling fee	1,00	65,00	65,00+	EUR	1,000000	0%	65,00+ (1)
Carrier charges	1,00	35,00	35,00+	EUR	1,000000	0%	35,00+ (1)
Insurance premium and tax	1,00	65,00	65,00+	EUR	1,000000	0%	65,00+ (1)
Insurance administration charges	1,00	25,00	25,00+	EUR	1,000000	0%	25,00+ (1)

Totals per currency: 725,00 EUR

10936349

725,00 EUR

* All our operations are subject to the Belgian Freight Forwarding Standard Trading Conditions 2005 (Annex M.B. 3486/2005 nr. 36362/07). Total will be sent free of charge upon first request. National and international trucking, organized by us, is subject to the application of the Belgian Law dated 489/1992 ensuring the international convention of the contract for the international carriage of goods by road (CMR). Our warehouse activities are governed by the General Logistic Conditions, dated October 9th, 2015 at the Office of the Chamber of Commerce and Industry of Antwerp and Brussels.

* In case of delay in the settlement of the invoice, the costs for recovery shall be due, as well as a yearly interest of 10% and an increase of 10% as a fixed compensation.

* Any legal proceedings shall be within the exclusive jurisdiction of the Antwerp Courts.

If you would like to receive your invoices by email in future, please contact our administration. Contact details in the right top corner of this invoice.

Type : Outgoing invoice
 Invoice : 027/908601

Company : 031
 Dept. : 610
 File : 1632411

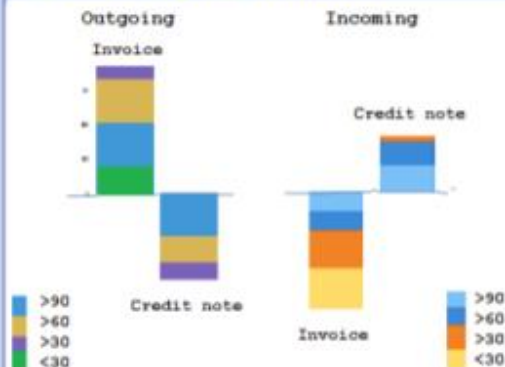
Customer: S16359
 STEPAN DEUTSCHLAND
 RODENKIRCHENER STR. 400
 50389 WESSELING
 DEU GERMANY

AMOUNT : 725,00 EURO

Related documents



Customer related invoices



Flow of functionality.

 = Clickable area

Document Preview

HERFURTH LOGISTICS

Document No: 007900001 Date: 01/10/2019

STEFAN SCHULZLAND GMBH
KOLN, 50669 KOLN, GERMANY

ORIGINAL

Document Details

Type: Transport document
Document: 007900001

Company: 400
Region: 110
Plant: 1000001

Customer: 1000001
Customer Name: STEFAN SCHULZLAND GMBH
Address: KOLN, 50669 KOLN, GERMANY

Related documents

Customer selected documents

Document Details

Document: 007900001 Date: 01/10/2019

Document Details

Type: Transport document
Document: 007900001

Company: 400
Region: 110
Plant: 1000001

Customer: 1000001
Customer Name: STEFAN SCHULZLAND GMBH
Address: KOLN, 50669 KOLN, GERMANY

Related documents

Customer selected documents

Document Details

Document: 007900001 Date: 01/10/2019

Index	Item	ItemName	ItemName	ItemName	ItemName
1	1000001-000001	Transport	Transport	Transport	Transport
2	1000001-000002	Transport	Transport	Transport	Transport
3	1000001-000003	Transport	Transport	Transport	Transport
4	1000001-000004	Transport	Transport	Transport	Transport
5	1000001-000005	Transport	Transport	Transport	Transport
6	1000001-000006	Transport	Transport	Transport	Transport
7	1000001-000007	Transport	Transport	Transport	Transport
8	1000001-000008	Transport	Transport	Transport	Transport
9	1000001-000009	Transport	Transport	Transport	Transport
10	1000001-000010	Transport	Transport	Transport	Transport
11	1000001-000011	Transport	Transport	Transport	Transport
12	1000001-000012	Transport	Transport	Transport	Transport
13	1000001-000013	Transport	Transport	Transport	Transport
14	1000001-000014	Transport	Transport	Transport	Transport
15	1000001-000015	Transport	Transport	Transport	Transport
16	1000001-000016	Transport	Transport	Transport	Transport
17	1000001-000017	Transport	Transport	Transport	Transport
18	1000001-000018	Transport	Transport	Transport	Transport
19	1000001-000019	Transport	Transport	Transport	Transport
20	1000001-000020	Transport	Transport	Transport	Transport

D.S.O. Customer(s)								
Customer C#	Customer Name	Netto Outstandings	Trend Outstand	DSO Factor	DSO Trend	In Limits	Expired < 30	Expired
1	A00504 AKOCHAR SHIPPING BV	2,526,58 EUR	57.37 %	30	-6.54 %	2,526,58 EUR	0,00 EUR	
2	A10797 MOL LOGISTICS (JAPAN) CO LTD	68,527,68 EUR	63.95 %	63.47	-24.74 %	52,617.78 EUR	7,909.90 EUR	
3	A11885 ALLTRANS FREIGHT & LOG	0,00 EUR	-100 %	0	-100 %	-100.42 EUR	306.42 EUR	
4	A12850 AGS COUSSAERT BELGIUM SPRL	1,695,00 EUR	-42.08 %	59.52	6.32 %	868.00 EUR	338.00 EUR	
5	A12890 AISEN EUROPE SA	495,00 EUR	48.25 %	14.82	-59.6 %	0,00 EUR	495.00 EUR	
6	A13237 ARTILAT BV	4,494,08 EUR	-78.96 %	50.33	67.77 %	2,678.58 EUR	1,815.50 EUR	
7	A12548 ADDV LTD	377,65 EUR	-6.15 %	30	0 %	0,00 EUR	377.65 EUR	
8	A14578 ASTRACON INTERNATIONAL	5,949,00 EUR	0 %	30	0 %	0,00 EUR	5,949.00 EUR	
9	A14622 AZC INTERNATIONAL INC.	668,00 EUR	-88.47 %	0	0 %	0,00 EUR	0,00 EUR	
10	A14739 ANGLO CARGO INTL LTD	3,590,05 EUR	108.71 %	19.65	0 %	0,00 EUR	5,482.13 EUR	
11	A14806 AFRICAN DESK	430,00 EUR	0 %	30	0 %	0,00 EUR	430.00 EUR	

Recognize text

Results



Words	Score
ibm	0.672
code	0.978
like	0.982
a	0.813
startup	0.989
ibm	0.947

Recognize tekst - Watson



```
"images": [
  {
    "image": "img_0235.jpg",
    "text": "62 [09]\n22",
    "words": [
      {
        "line_number": 0,
        "location": {
          "height": 48,
          "left": 322,
          "top": 227,
          "width": 49
        },
        "score": 0.5176,
        "word": "62"
      },
      {
        "line_number": 0,
```

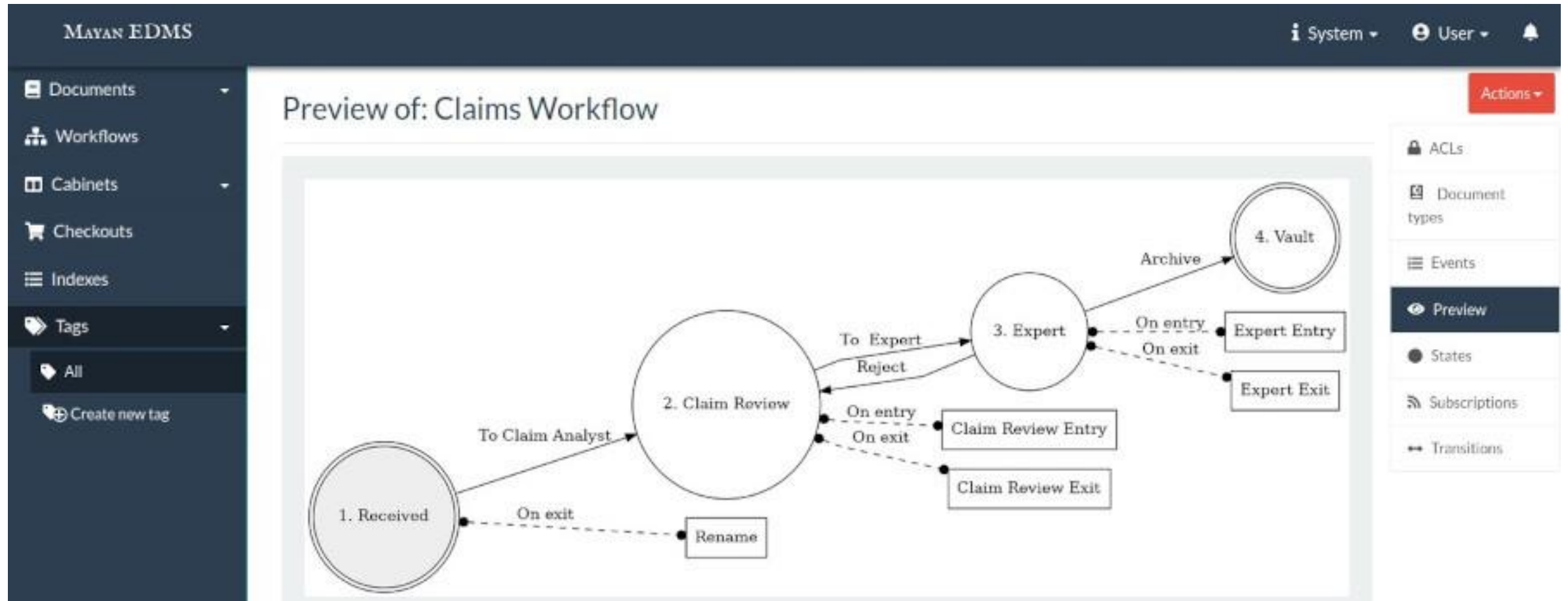

Recognize tekst – Tesseract



```
"textAnnotations": [  
  {  
    "locale": "en",  
    "description": "YMLU 621093 9\n22U1\nCAUTION\n07/08/2017",  
    "boundingPoly": {  
      "vertices": [  

```

Mayan EDMS



Mayan EDMS

The screenshot displays the Mayan EDMS web application interface. The top header bar is dark blue with the text 'MAYAN EDMS' on the left and 'System', 'User', and a notification bell icon on the right. A left sidebar contains a 'Documents' menu with options like 'Recently accessed', 'Recently added', 'Favorites', 'All documents', 'Trash can', 'Duplicated documents', 'New document', 'Workflows', 'Cabinets', 'Checkouts', 'Indexes', and 'Tags'. The main content area shows the title 'OCR result for document: CMS1500ClaimForm010402.pdf' and a large text box containing the OCR text. A right sidebar lists various actions: 'ACLs', 'Cabinets', 'Check in/out', 'Comments', 'Content', 'Duplicates', 'Events', 'File metadata', 'Indexes', 'Metadata', 'OCR' (highlighted), 'Pages', and 'Preview'. A red 'Actions' button is located at the top right of the main content area.

MAYAN EDMS

System User

Documents

- Recently accessed
- Recently added
- Favorites
- All documents
- Trash can
- Duplicated documents
- New document
- Workflows
- Cabinets
- Checkouts
- Indexes
- Tags

OCR result for document: CMS1500ClaimForm010402.pdf

Use a participating provider to receive the maximum benefit. Durable medical equipment and ongoing services as physical therapy are especially cost effective with a UnitedHealthcare provider.

Please review your benefits at myuhc.com. For services that require prior authorization or notification, be the Member Services number on the back of your health plan ID card.

What happens next:

After we process your claim, we will send you an Explanation of Benefits (EOB). The EOB will explain the c applied to your plan deductible and any charges you owe your health care provider. Please keep your EOB on future reference. You also may review your EOB information online at myuhc.com.

Once you have completed the form, mail it to the address listed on the back of your Health Plan ID Card. Be sure to attach the Superbill or Invoice and any receipts of your payments.

Actions

- ACLs
- Cabinets
- Check in/out
- Comments
- Content
- Duplicates
- Events
- File metadata
- Indexes
- Metadata
- OCR
- Pages
- Preview

Mayan Full Text Search

The screenshot displays the Mayan EDMS interface. On the left is a dark sidebar with navigation links: Documents, Workflows, Cabinets, Checkouts, Indexes, and Tags. The Tags section is expanded, showing 'All' and 'Create new tag'. The main area is titled 'Dashboard' and features a search bar with the text 'what is this form'. Below the search bar are eight summary cards arranged in a 3x3 grid (with the last cell empty). Each card includes an icon, a title, a value, and a 'View details' link with an external link icon.

Icon	Category	Value
	Checked out documents	0
	Total documents	4
	Total pages	28
	Documents in trash	3
	Document types	3
	New documents this month	1
	New pages this month	11
	Total groups	6
	Total roles	5
	Total users	9

Tesseract best practices

- Provide input file with high quality

If possible, provide the software with a high quality input file. Poor image or document quality may prevent Tesseract from recognizing the text correctly. This is also true when processing documents with complex structures. Tesseract has problems recognizing complex structures such as tables and mixed text-image documents.

- Perform preprocessing

Perform appropriate preprocessing of the image data, such as contrast adjustment, noise reduction, and sharpening, to improve Tesseract's text recognition performance.

Tesseract best practices

- Define Region of Interest (ROI)

Define a region of interest around the relevant text area to increase recognition accuracy and reduce processing time.

- Make language selection

Make sure that the language setting of Tesseract matches the detected language in the image for best results.

Tesseract best practices

- Perform model training

If needed, you can improve Tesseract OCR by training a custom model for specific text types or fonts. This enables more accurate text recognition in specific scenarios.

- Perform validation and error correction

Review and correct recognized text results. Use validation tools and implementations for automatic error correction to improve the quality of recognized texts.

Tesseract resources

- The official documentation for Tesseract
- <https://tesseract-ocr.github.io/>
- <https://github.com/invoice-x/invoice2data>
- <https://regex101.com/>

Conclusion

Conclusion

After having fun with Tesseract OCR, I can say that the engine is amazing!! Here the list of interesting point from Tesseract in my opinion:

- Cost efficient - no expensive licenses
- Open Source.
- Easy to use.
- Good extract result.
- Support multi language (Latin & Non-Latin).