

,

()

()

()

,

()

:

II

,

-20-1

(,)

123 «

'

»

()

-

(- -)

,

()

:

(, ,)

()

(,)

,

()

123 « ' »

()

-

(- -)

,

()

:

“ ” 20 .

(, ,)

1.

,

“ 05 ” 2021 . 1656

2.

13

2020 .

3.

,

,

4.

,

5. _____ , _____ , _____ , _____ , _____
 () 11 _____

6. _____ , _____ .1) (_____)

	(_____ , _____ , _____ , _____)		

1		09.11.21-14.11.21	
2		15.11.21-21.11.21	
3		22.11.21-30.11.21	
4		01.12.21-06.12.21	
5		07.12.21-09.12.21	
6		09.12.21-10.12.21	

8 2021 .

_____ ()
 | _____ () _____ (, ,)

: 70 ., 31 ., 6 .,

1 ., 13 .

OCR,

, ,
OPENCV, TESSERACT, SPACY.

ABSTRACT

Master's thesis: 70 pages, 31 figures, 6 tables, 1 appendices, 13 sources.

OCR, IMAGE PROCESSING, NOISE REDUCTION, TEXT ANALISYS,
LEMMATIZATION, STEMMING OPENCV, TESSERACT, SPACY

The major goal of this thesis is is to study the impact of the use of different methods of segmentation and noise reduction on the accuracy and speed of recognition of scanned documents.

In order to the qualification work, the existing methods of image processing for noise reduction and segmentation were considered, the analysis of the impact of these methods on the accuracy of the text recognition system was analyzed.

A system of orderly storage of scanned documents with keyword selection was also developed, which consists of an image processing module, a text recognition module and modules for text analysis. In addition, image and text processing methods have been adapted for different types of computers.

	,	,	,	
			8
			9
1			10
1.1			11
1.2			13
1.2.1	ABBY Finereader.....			13
1.2.2	ABBY Flexicapture.....			14
1.2.3	Adobe Acrobat DC.....			14
1.3			15
2				
			17
2.1			17
2.1.1			18
2.1.2			19
2.1.3			20
2.2			21
2.2.1			21
2.2.2			22
2.3			23
2.3.1			24
2.3.2			25
2.3.3			26
2.3.4			26
2.4			27
2.4.1	OpenCV.....			27
2.4.2	Tesseract.....			30

2.5	32
2.6	32
2.7	33
2.8	CUDA.....	34
2.9	38
3		
	40
3.1		
	40
3.2	44
3.3	47
3.3.1	48
3.3.2	48
3.3.3 POS	49
3.3.4	,	52
3.3.5	-	52
3.3.6 TF-IDF	53
3.3.7	54
4		
	56
4.1	56
4.2	58
4.3		
	59
	61
	62
	64

‘ ‘ ‘

OCR –

POS –

(part-of-speech)

LSTM –

’ (Long short-term memory)

[2].

1.1

(OCR) –

OCR –

OCR

OCR

OCR

()

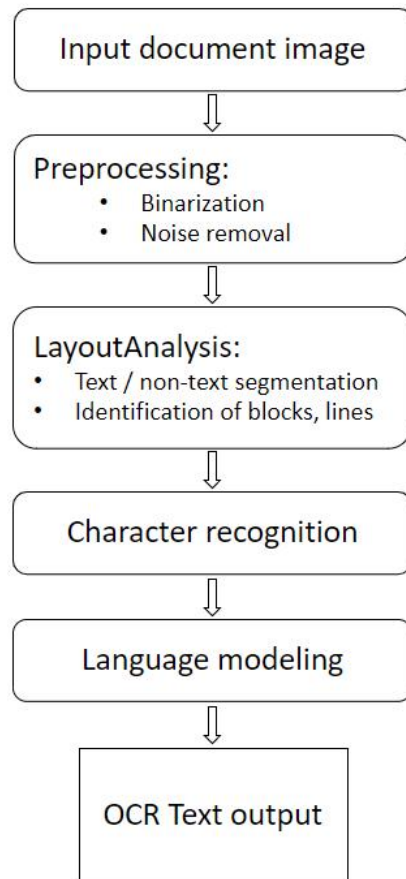
[3].

OCR

OCR

[4].

1.1



1.2

,

OCR

/

,

.

,

-

,

,

,

OCR

,

/

.

,

,

[5].

1.2.1 ABBYY Finereader

ABBYY FineReader PDF –

PDF.

Finereader:

- ;
- ;
- .

Finereader:

- ;
- ;
- .

1.2.2 ABBYY Flexicapture

FlexiCapture –

- ,
- ,
- Flexicapture :
- ;
- ;
- ().
- Flexicapture :
- ;
- /
- ;
- ;
- .

1.2.3 Adobe Acrobat DC

Adobe Acrobat DC –

- PDF- , PDF-
- OCR. ,
- ,
- .
- , Adobe Acrobat
- ,
- :
- ;
- .

2

(OCR).

OCR

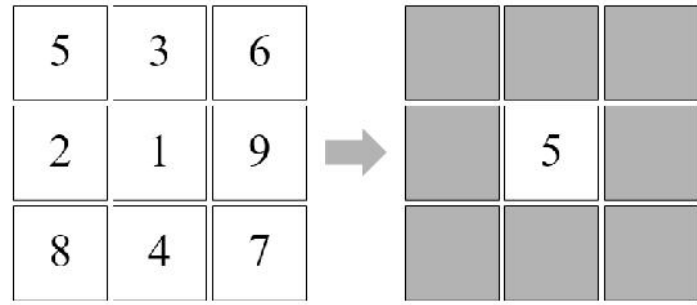
OCR

[2].

2.1

2.1.1

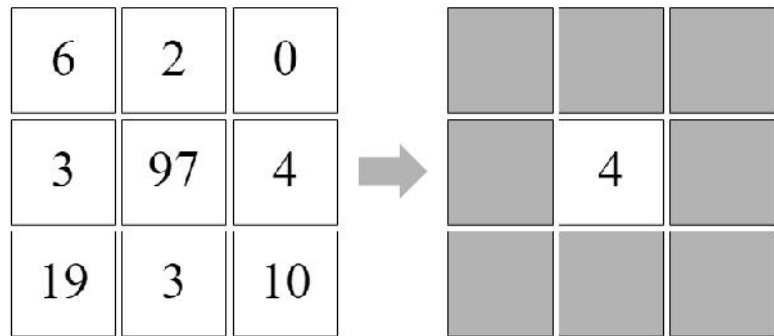
3x3



2.1 –

2.1.2

– ,
 . ,
 .
 .
 , ()



2.2 –

, (2.1)
 5,
 1, 2, 3, 4, 5, 6, 7, 8, 9.
 16, , 144 144/9 = 16.
 :

», «

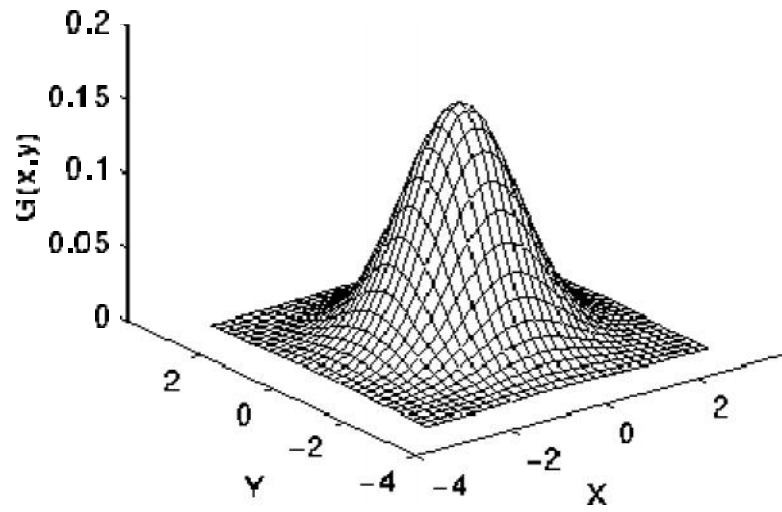
2.1.3

, , (x y), (2.1).

$$G(x, y) = \frac{1}{2f\uparrow^2} e^{-\frac{x^2+y^2}{2\uparrow^2}}. \quad (2.1)$$

()

(0,0) = 1 2.3.



2.3 –

2.2

0,

(255).

2.2.1

2.2.2

... ,
 ... ,

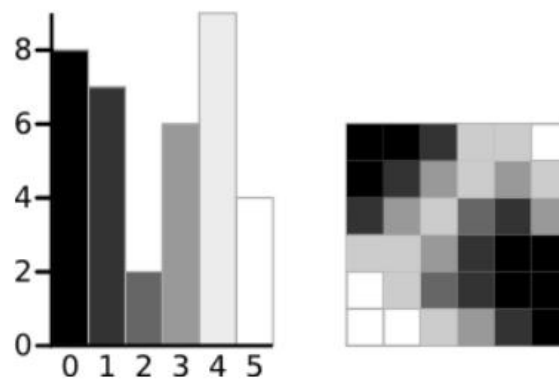
 ... ,
 ... ,
 ... (...)
 ... ,
 ... [10].

t :

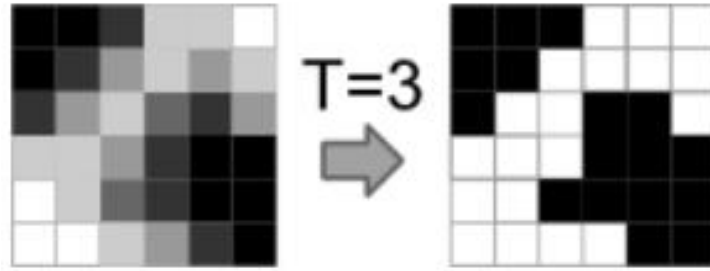
$$\dagger^2(t) = \check{S}_{bg}(t) \dagger_{bg}^2(t) + \check{S}_{fg}(t) \dagger_{fg}^2(t), \quad (2.2)$$

$\check{S}_{bg}(t)$ $\check{S}_{fg}(t)$ - ;

\dagger^2 - .



2.4 -

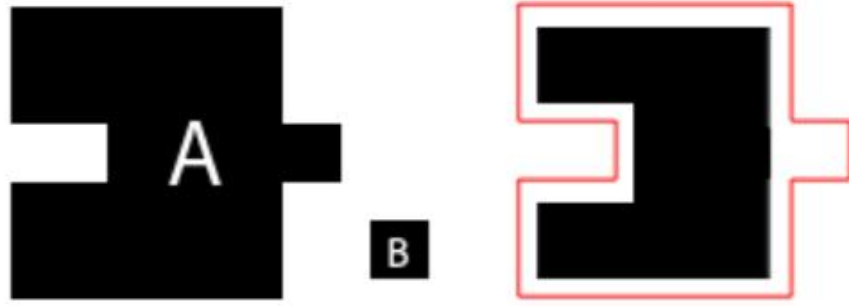


2.5 -

2.3

[11].

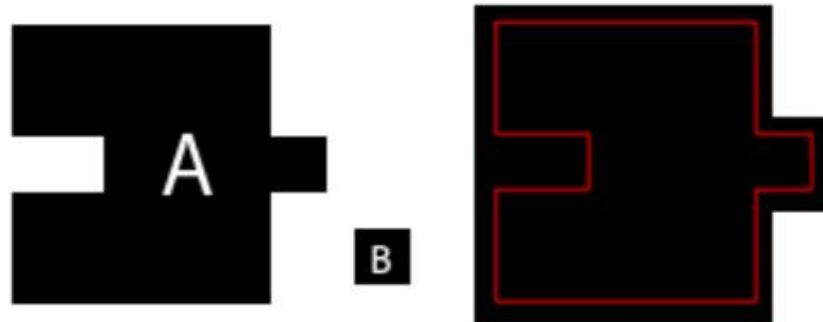
« »



2.6 -

2.3.2

. "1",
 "1". ,
 ,
 . ,
 . ,
 , .. ,
 . ,
 . ,
 .



2.7 -

A

B

:

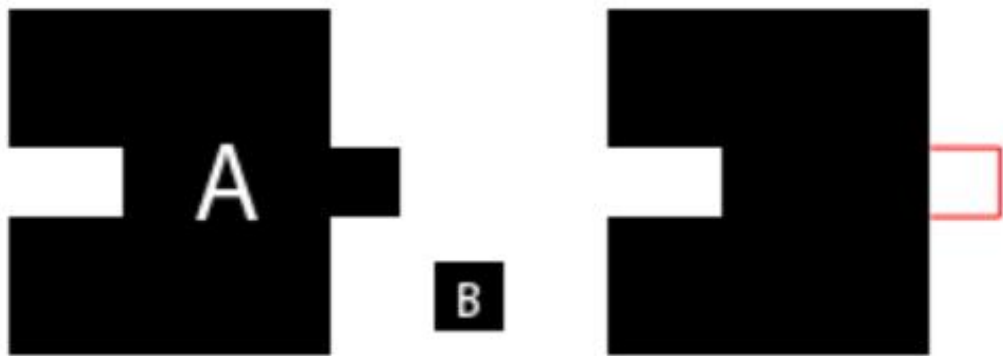
$$A \oplus B \{z \in E \mid (B^s)_z \cap A \neq \emptyset\}$$

(2.4)

2.3.3

,
 , - ,
 , -
 . A B A B
 B.
 A B

$$A \circ B = (A \ominus B) \oplus B. \tag{2.5}$$

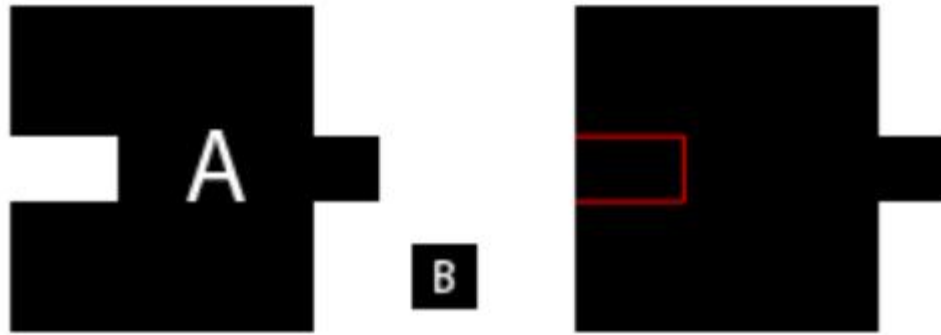


2.8 -

2.3.4

,
 ,
 . A B
 A B, B.
 A B

$$A \bullet B = (A \oplus B) \ominus B. \tag{2.6}$$



2.9 –

2.4

2.4.1 OpenCV

OpenCV (Open Source Computer Vision Library) –
(API), Intel,

. OpenCV 1999 ,
Intel Research

,

3D-

OpenCV –

. OpenCV

. OpenCV

C

OpenCV 500 ,

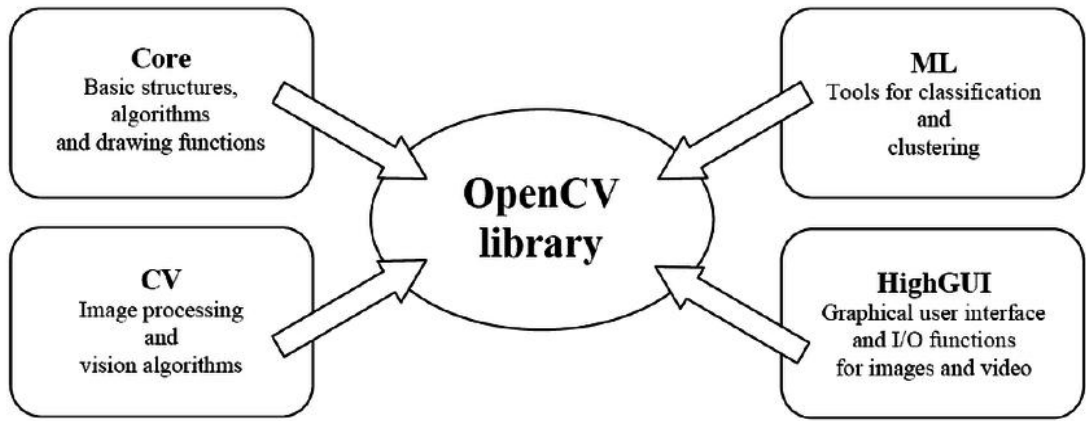
, , , , ,

. ,
 ,
 « - ' », , , ,
 . OpenCV
 ,
 . OpenCV DirectX,
 API, Microsoft
 .
 OpenCV .
 Borland C++, Microsoft Visual Studio C++
 Intel. C C++ ,
 . OpenCV
 Windows, Linux. OpenCV
 ,
 . ,
 , , . ,
 ,
 . ,
 . -
 , ' - . OpenCV
 , ' . OpenCV
 , ,
 . ,
 , -
 , .

```

,
,
.
,
.
.
.
.
.
:
:
- ;
- ;
- .
, ' - .
« »
, , , , , ,
, , , , , ,
, , , , , ,
[7].
OpenCV ( 2.10):
- , (API
);
- , ;
- ;
- ;
- , ,
.

```



2.10 –

OpenCV

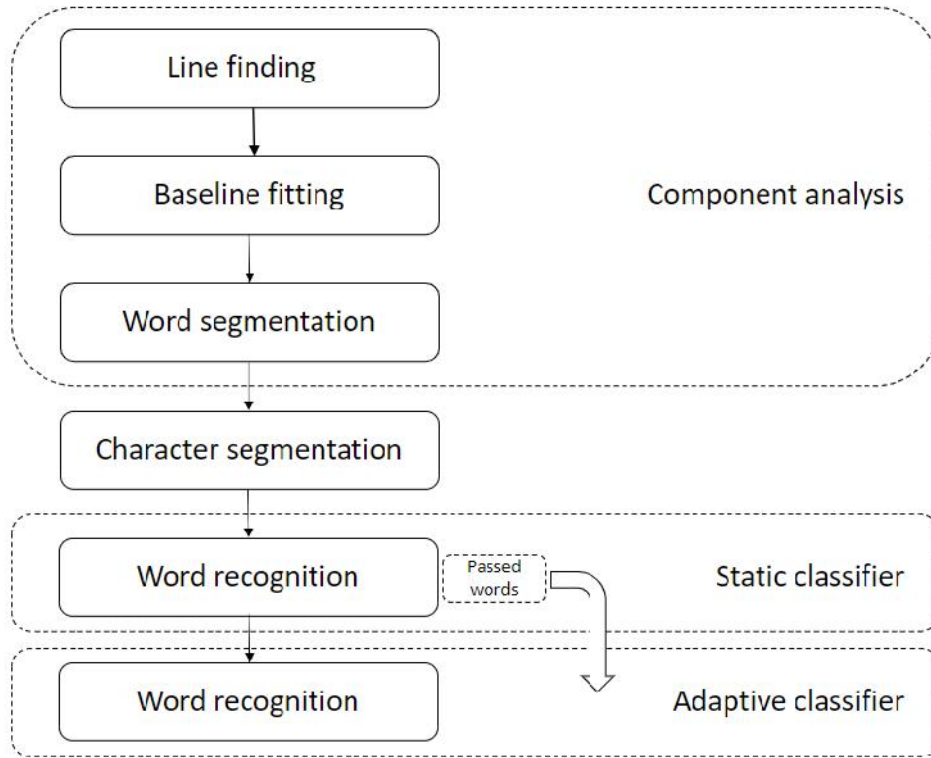
2.4.2 Tesseract

Tesseract – OCR
 HP 1984 1994 . Tesseract
 HPLabs,
 OCR

. Tesseract

. Tesseract

(Blob).



2.11 –

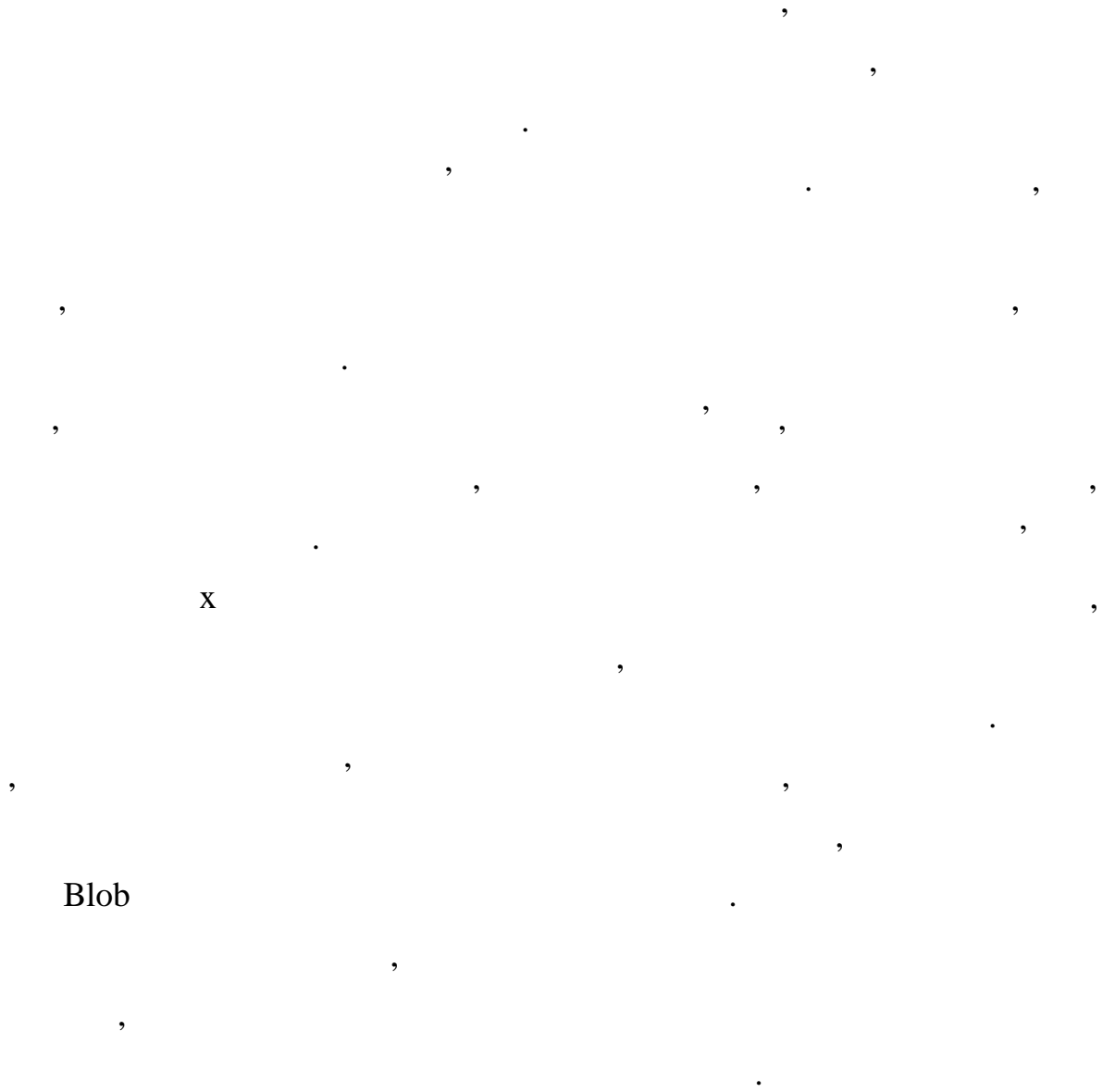
Tesseract

[8].

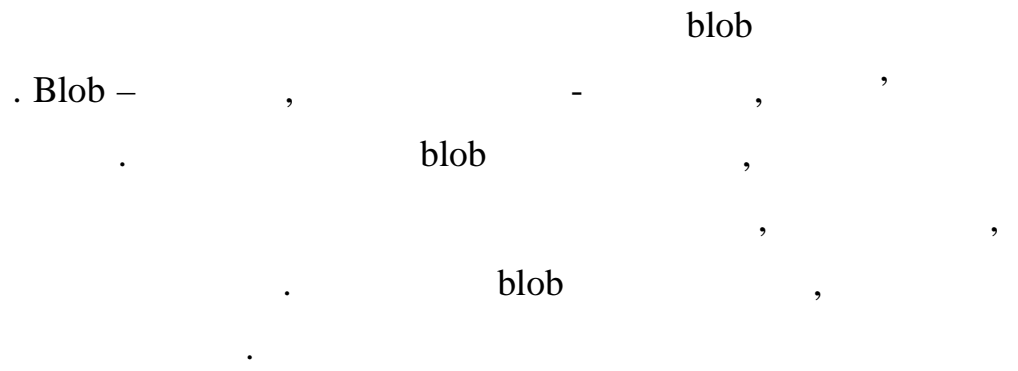
Tesseract

2.11.

2.5



2.6



blob. ; Tesseract
(2.12) [9].

Volume 69, pages 872-879.

2.12 -

2.7

, (),

. Tesseract

GPU -

(SM)

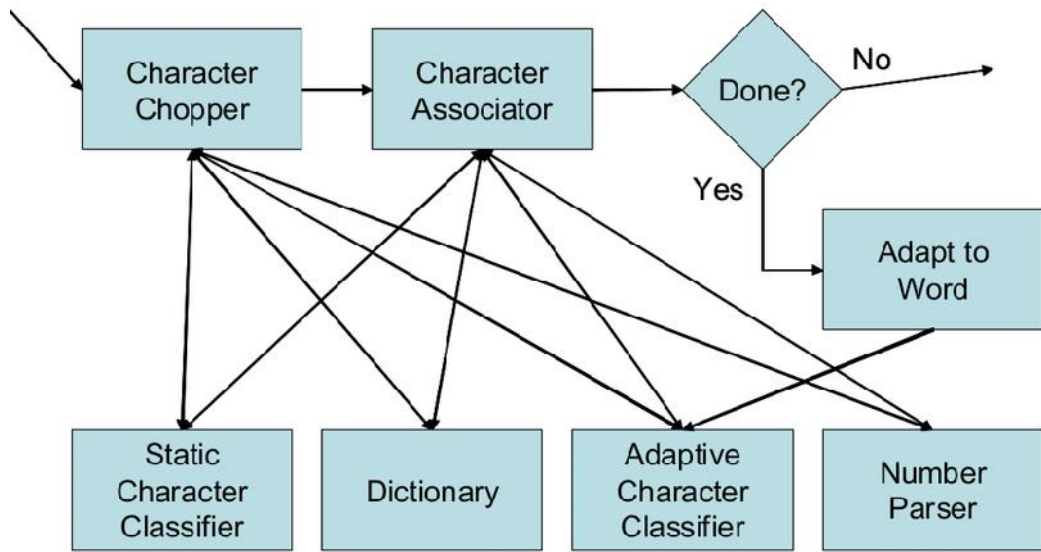
(SIMD). SM

SM

().

2.14.

HD (),



2.13 –

2.8

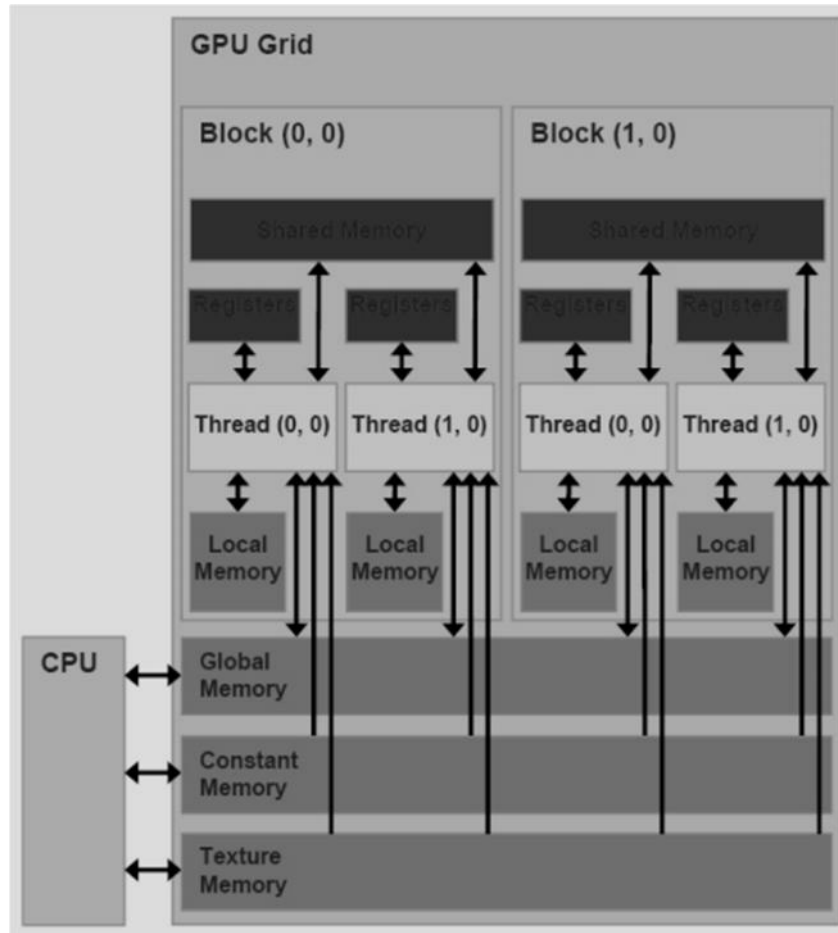
CUDA

CPU

GPU

CPU

GPU

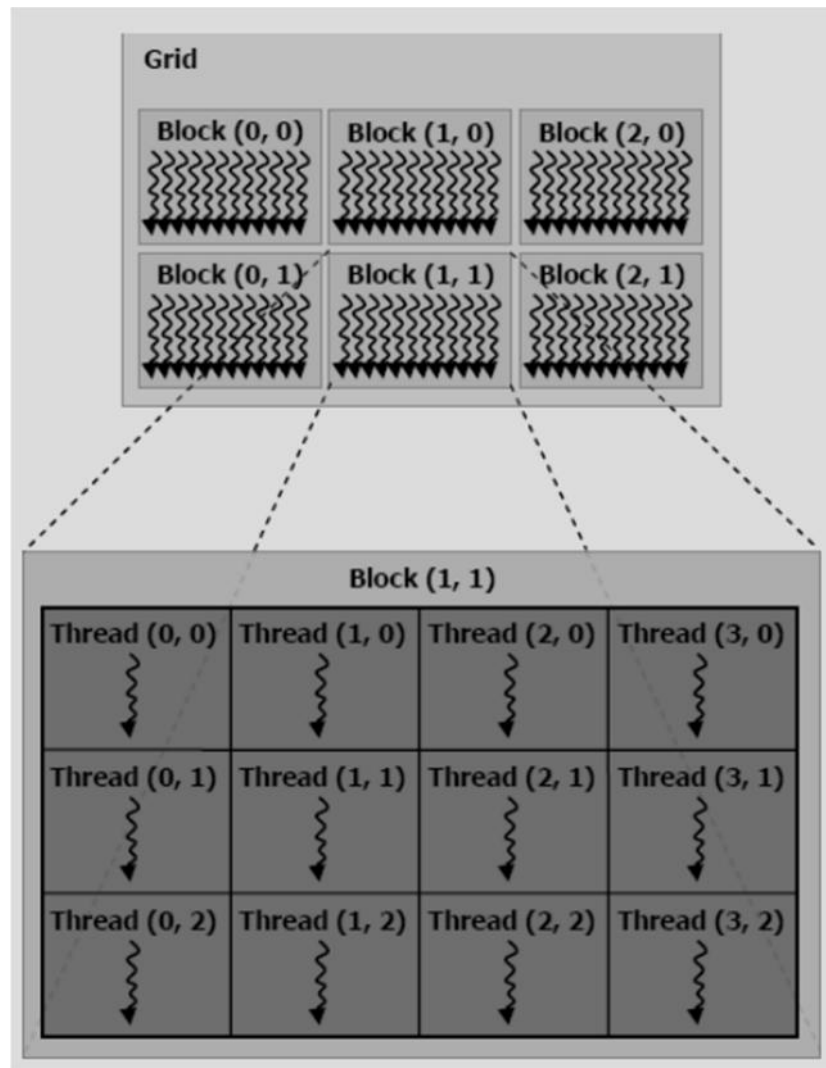


2.14 –

CUDA

GPU.

()



2.15 –

GPU

CUDA

, (2.16).

, , . SM

, , . ;

, . , -

, , . ,

, . ,

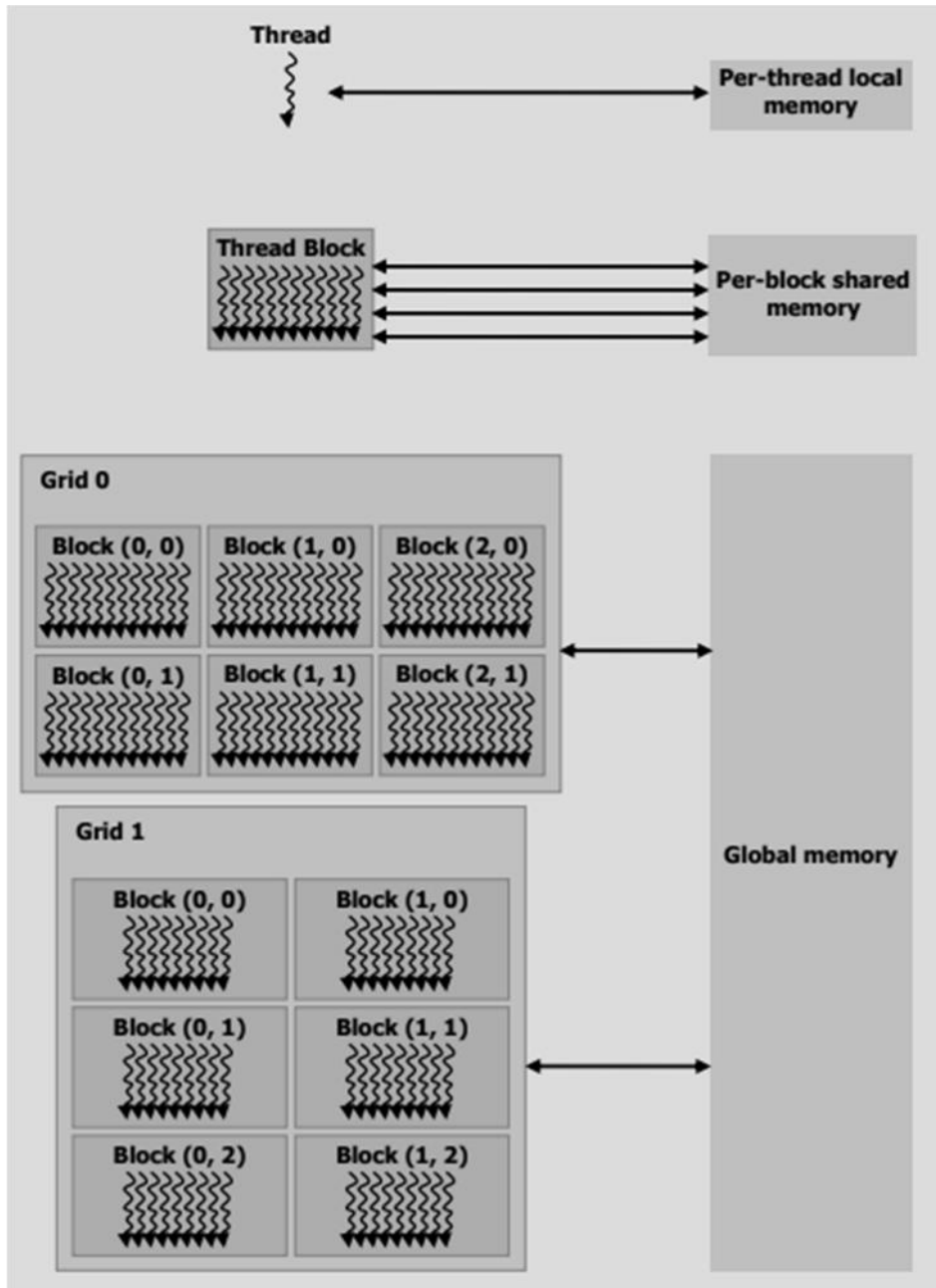
, , . ,

, . ,

, . ,

, . ,

, . ,



2.16 –

CUDA

2.9

(POS)

(NLP).

(POS)

. POS-
 ,
 . , POS
 .
 POS-
 . POS-
 (TTS),
 , , ,
 POS , -
 , ,
 , POS
 , ,
 POS ,
 , ,
 , ,
 , ,
 . ,
 ,

3

3.1

.
 ,
 .
 , , , ,
 .
 ,
 .
 :
 - ;
 - ;
 - .
 ,
 ,
 :
 - ;
 - ;
 - ;
 - .

,

—

OCR

,

OCR,

OCR

OCR.

:

- ;

- Pos ;

- - ;

- ;

- .

,

.

.

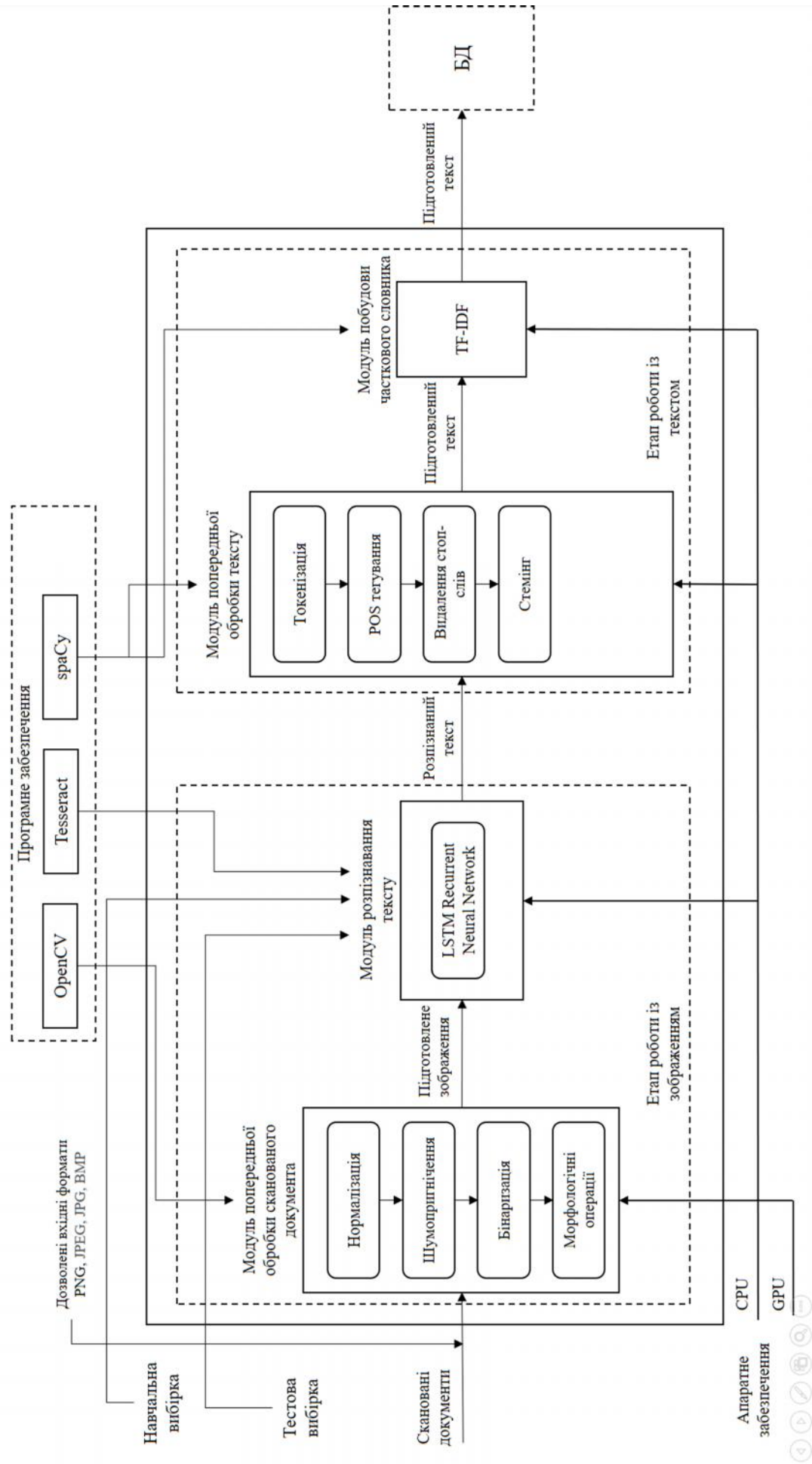
,

saw,

s,

see,

saw



3.1 –

3.2

(4.4)

4.5 4.6

In 1830 there were but twenty-three miles of railroad in operation in the United States, and in that year Kentucky took the initial step in the work west of the Alleghanies. An Act to incorporate the Lexington & Ohio Railway Company was approved by Gov. Metcalf, January 27, 1830. It provided for the construction and re-

4.4 -

In 1830 there were but twenty-three miles of railroad in operation in the United States, and in that year Kentucky took the initial step in the work west of the Alleghanies. An Act to incorporate the Lexington & Ohio Railway Company was approved by Gov. Metcalf, January 27, 1830. It provided for the construction and re-

4.5 -

In 1830 there were but twenty-three miles of railroad in operation in the United States, and in that year Kentucky took the initial step in the work west of the Alleghanies. An Act to incorporate the Lexington & Ohio Railway Company was approved by Gov. Metcalf, January 27, 1830. It provided for the construction and re-

4.6 -

(4.7).

In 1830 there were but twenty-three miles of railroad in operation in the United States, and in that year Kentucky took the initial step in the work west of the Alleghanies. An Act to incorporate the Lexington & Ohio Railway Company was approved by Gov. Metcalf, January 27, 1830. It provided for the construction and re-

4.7 -

(4.6).

In 1830 there were but twenty-three miles of railroad in operation in the United States, and in that year Kentucky took the initial step in the work west of the Alleghanies. An Act to incorporate the Lexington & Ohio Railway Company was approved by Gov. Metcalf, January 27, 1830. It provided for the construction and re-

4.8 -

In 1830 there were but twenty-three miles of railroad in operation in the United States and in that year Kentucky took the initial step in the work west of the Alleghanies. An Act to incorporate the Lexington & Ohio Railway Company was approved by Gov. Metcalf, January 27, 1830. It provided for the construction and re-

4.9 -

3.3

spaCy.

(NLP) Python.

spaCy

spaCy

spaCy

Python.

- , ,

- , -

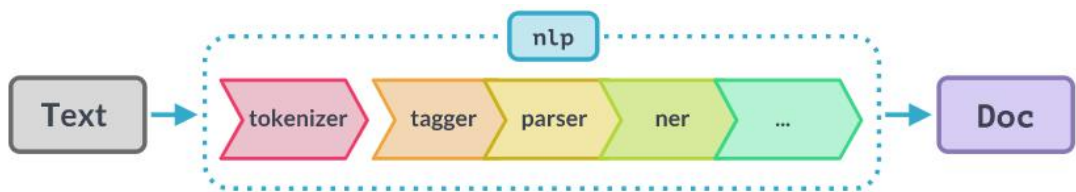
, , spaCy

3.3.1

, .

. spaCy

,



3.7

3.3.2

: «"Elon Musk built his electric car company, Tesla, around the promise that it represented the future of driving".

Elon
Musk
built
his
electric
car
company
,
Tesla
,
around
the
promise
that
it
represented
the
future
of
driving

3.8 –

3.3.3 POS

«the»

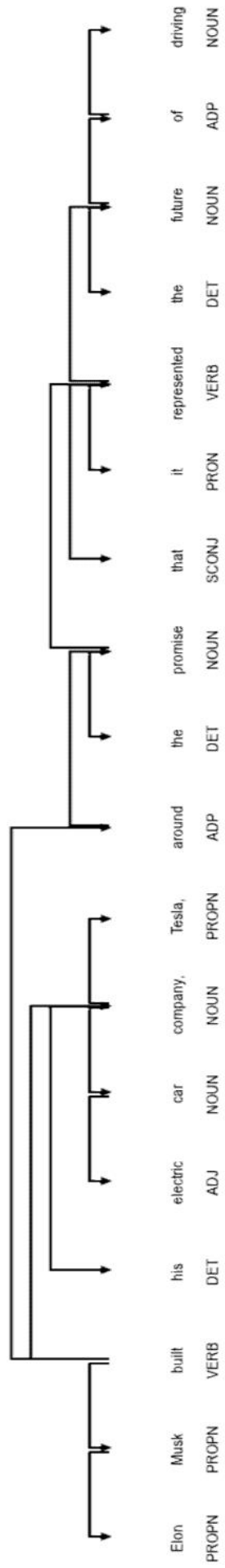
Token.

```

:
- Text: .
- Lemma: .
- POS: UPOS.
- Tag: .
- Dep: , .

```

TEXT	LEMMA	POS	TAG	DEP
Elon	Elon	PROPN	NNP	compound
Musk	Musk	PROPN	NNP	nsubj
built	build	VERB	VBD	ROOT
his	-PRON-	DET	PRP\$	poss
electric		electric	ADJ	JJ amod
car	car	NOUN	NN	compound
company	company		NOUN	NN dobj
,	,	PUNCT	,	punct
Tesla	Tesla	PROPN	NNP	appos
,	,	PUNCT	,	punct
around	around	ADP	IN	prep
the	the	DET	DT	det
promise	promise		NOUN	NN pobj
that	that	SCONJ	IN	mark
it	-PRON-	PRON	PRP	nsubj
represented		represent	VERB	VBD acl
the	the	DET	DT	det
future	future	NOUN	NN	dobj
of	of	ADP	IN	prep
driving	driving		NOUN	NN pcomp



3.1 –

```
from spacy.lang.en.stop_words import STOP_WORDS

for word in token_list:
    lexeme = nlp.vocab[word]
    if lexeme.is_stop == False:
        filtered_sentence.append(word)
```

3.2 –

Source

```
['Elon', 'Musk', 'built', 'his', 'electric', 'car', 'company',
 'Tesla', 'around', 'the', 'promise', 'that', 'it',
 'represented', 'the', 'future', 'of', 'driving']
```

Without stop words

```
['Elon', 'Musk', 'built', 'electric', 'car', 'company', 'Tesla',
 'promise', 'represented', 'future', 'driving']
```

3.3.6 TF-IDF

3.3 –

```
def get_frequency_matrix(sentences):
    frequency_matrix = {}
    ps = PorterStemmer()
    stopWords = set(stopwords.words("english"))

    for sent in sentences:
```

```

table = {}
words = nltk.word_tokenize(sent)

for word in words:
    word = word.lower()
    word = ps.stem(word)
if word in stopWords:
    continue

if word in table:
    table[word] += 1
else:
    table[word] = 1

frequency_matrix[sent[:15]] = table

return frequency_matrix

```

3.4–

```

{'Elon': {'elon': 1},
'Musk': {'musk': 1},
'built': {'built': 1},
'electric': {'electr': 1},
'car': {'car': 1},
'company': {'compani': 1},
'Tesla': {'tesla': 1},
'promise': {'promis': 1},
'represented': {'repres': 1},
'future': {'futur': 1},
'driving': {'drive': 1}}

```

3.3.7

```

,
.
,
Python.
nlp.pipe.
,
Doc.
Doc,
.

```


4

4.1.

4.1 –

CPU	Intel Core i5-4210U (1.7 - 2.7)
GPU	nVidia GeForce 840M
GPU ()	384 CUDA, 1029 , 2 DDR3
RAM	12 Gb
OS	Windows 10

$$Speedup = \frac{T_s}{T_p} \tag{4.1}$$

T_s – ;

T_p – .

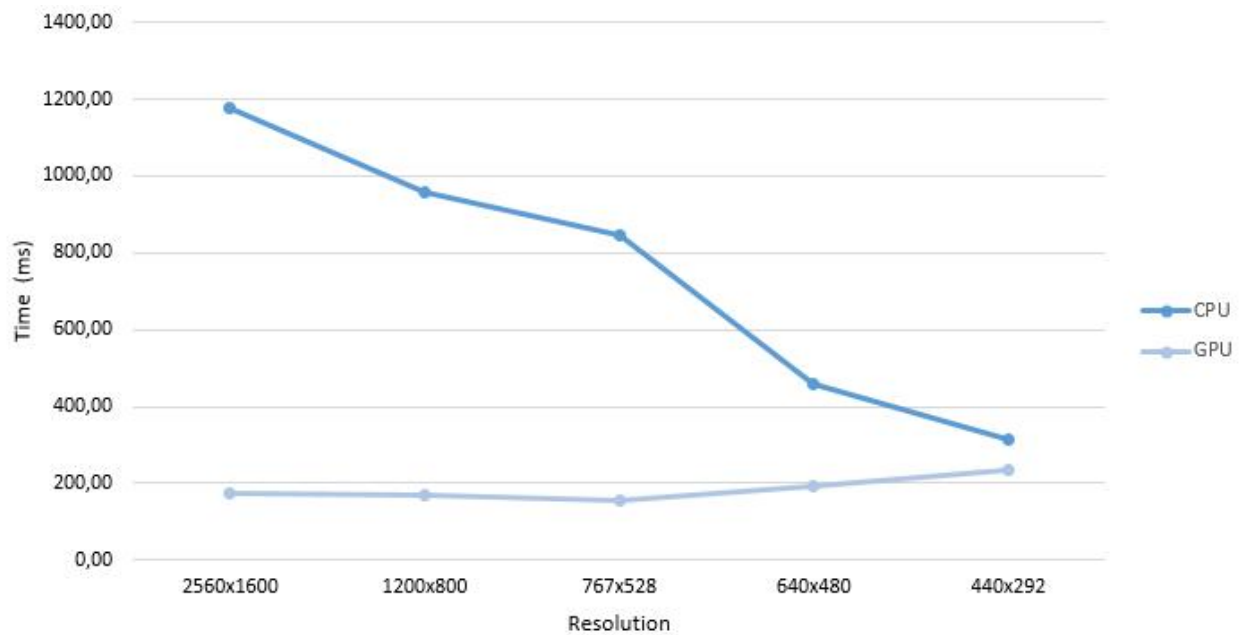
4.1

GPU.

4.2.

4.2 –

	CPU	GPU	
2560 1600	1176,00	173,45	6,780051888
1200 800	955,33	167,67	5,697813121
767 528	844,00	153,76	5,489073881
640 480	456,67	190,43	2,398081535
440 292	315,33	234,67	1,34375



4.1 –

(4.1)

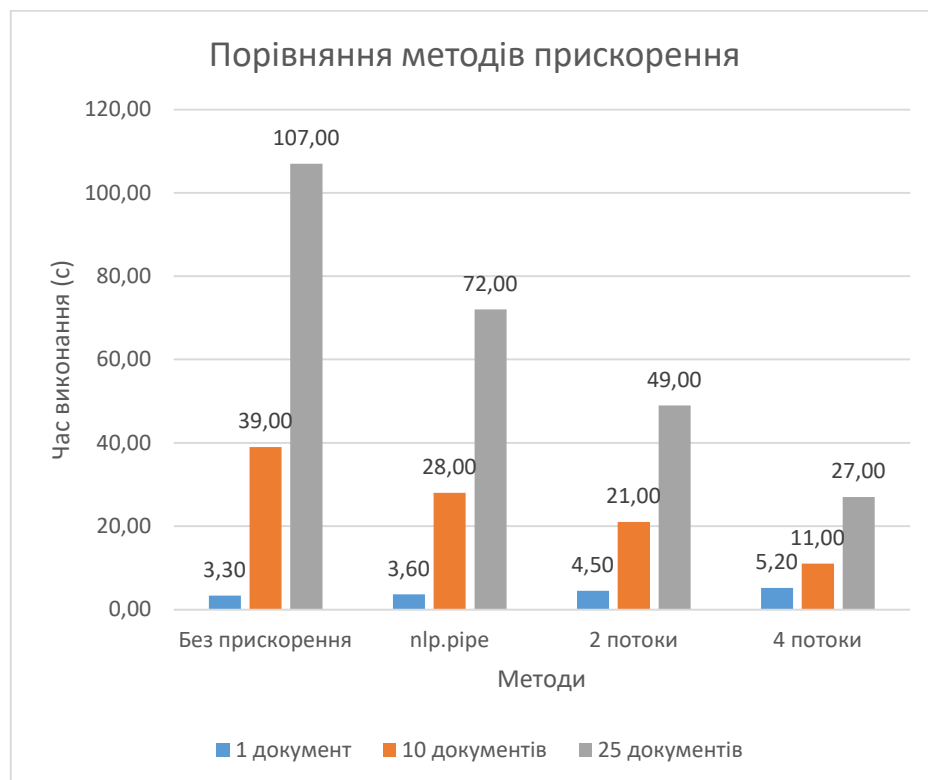
4.2

nlp.pipe

4.4.

4.4 –

		1	10	25
		3.5	39	107
nlp.pipe		3.6	28	72
	2	4.5	21	49
	4	5.2	11	27



4.2 –

(4.2),

nlp.pipe.

4.3

4.2.

4.2 –

	Pipeline			
	en_core_web_sm	en_core_web_md	en_core_web_lg	en_core_web_trf
	Vocabulary syntax entities	Vocabulary syntax tax entities vectors	Vocabulary syntax entities vectors	Vocabulary syntax tax entities
	13 MB	43 MB	741 MB	438 MB
	0	20	685	0

4.3,

en_core_web_sm,

en_core_web_trf.

4.3 –

	pipeline			
	en_core_w eb_sm	en_core_w eb_md	en_core_w eb_lg	en_core_we b_trf
	1.00	1.00	1.00	1.00
Pos	0.97	0.97	0.97	0.98
	0.92	0.91	0.92	0.96
	0.9	0.9	0.9	0.94
,	0.85	0.85	0.86	0.9

1. [] / . . . , . . . , – 2010. – : <https://core.ac.uk/download/pdf/74268943.pdf>.
2. Harraj A. E. OCR ACCURACY IMPROVEMENT ON DOCUMENT IMAGES THROUGH A NOVEL PRE-PROCESSING APPROACH / A. E. Harraj, N. Raissouni. – 2015.
3. Optical Character Recognition (OCR) [] . – 2019. – : <https://ukdiss.com/examples/optical-character-recognition.php>.
4. Optical Character Recognition – A Combined ANN/HMM Approach : . / Sheikh Faisal Rashid, 2014. – 161 .
5. Best OCR Software of 2021 [] . – 2021. – : <https://nanonets.com/blog/ocr-software-best-ocr-software/>.
6. Shrestha P. OPTICAL CHARACTER RECOGNITION [] / Prakash Shrestha. – 2018. – : https://www.theseus.fi/bitstream/handle/10024/151264/Shrestha_Pramoj.pdf?sequence=1&isAllowed=y.
7. Mumtazimah M. A Review on OpenCV [] / M. Mumtazimah, M. Yazid, S. H. Muhammad. – 2015. – : https://www.researchgate.net/publication/280977983_A_Review_on_OpenCV.
8. Ohlsson V. Optical Character and Symbol Recognition using Tesseract [] / Victor Ohlsson. – 2016. – : <http://www.diva-portal.org/smash/get/diva2:1019846/FULLTEXT02.pdf>.
9. Smith R. An Overview of the Tesseract OCR Engine [] / Ray Smith – : <https://static.googleusercontent.com/media/research.google.com/en//pubs/archive/33418.pdf>.

10. Otsu Thresholding [] – 2010. –
: <http://www.labbookpages.co.uk/software/imgProc/otsuThreshold.html>.
11. Morphological Image Processing [] – 2000. –
: <https://www.cs.auckland.ac.nz/courses/compsci773s1c/lectures/ImageProcessing-html/topic4.htm>.
12. Morphological Transformations [] – 2021. –
: https://docs.opencv.org/3.4.15/db/df6/tutorial_erosion_dilatation.html.
13. Morphology Filters [] –
: http://www.theobjects.com/dragonfly/dfhelp/4-0/Content/05_Image%20Processing/Morphology%20Filters.htm.