9~ 2010

(4)

(23)

_

جي

2009 / 02 / 16

2008 / 05 / 19

ABSTRACT

In this research we propose a statistical method and morpho-lexical analysis for correcting Arabic words as a post processor for Arabic words output from OCR systems. Dictionaries of words were built for the comparison to the attached word.

The present research uses multiple knowledge sources and basing on the Arabic language properties, statistical method, morpho-lexical analysis and dictionary look-up for error detection and correction. Correction of errors in this research depends on the type of possible error, which can be: transposing two adjacent letters, rejection, replacing an incorrect letter, inserting a missing letter, substitution errors, which are most frequently committed by the OCR systems.

.OCR

.1

.[1] :

(OCR- Optical Character Recognition)

C++

.2

: (n-gram)

n-gram

& & OCR [3,2] [4] [5] [6] .Viterbi [7] OCR (Morpho-Lexical) Production Rules [8] OCR .3 28 .[9] [10] (1 (2

(3

(4

.[11 9] (1 -)

:(1)

	.(1)	

& &

:[12] '(') (1 (2 .1 () (1 (2

.4

) Text

(...

. (2-)

:(2)

1372	4808	5361	3534	9819 ()	49514	
92	3204	1449	577	1081	6795 ()	
109	901	475 ()	155	673	4889	
213	565	637	517	1003	5441	
281	949	907 ()	447	1267	5079	
0	423	98 ()	82 ()	195 ()	1522	

.5

_ " "

_

и и

.[1]

& &

:[13] .6 .[1] (n) n- n = 1 4 3 2 1 n .(n-gram) diagram n=2 ; $mono\ grams$ uni-grams grams. trigrams n=3 ; bigramsn-grams n-grams n-grams .[12] .(...

136

.(...

.()	:	
		Notepad			
		•••			
:					
,				-	-6
(strtok))			\ C	(1
. (() ;)) C++	
	() , ,				(2
)	((3
	((4
					`
				•	
	•				(5
ASCII Code					•
				_	-6
		-:			
	•				1

& & .

: 1-

: 2-

:_______.2

. (1-)

÷

.(2)

.3

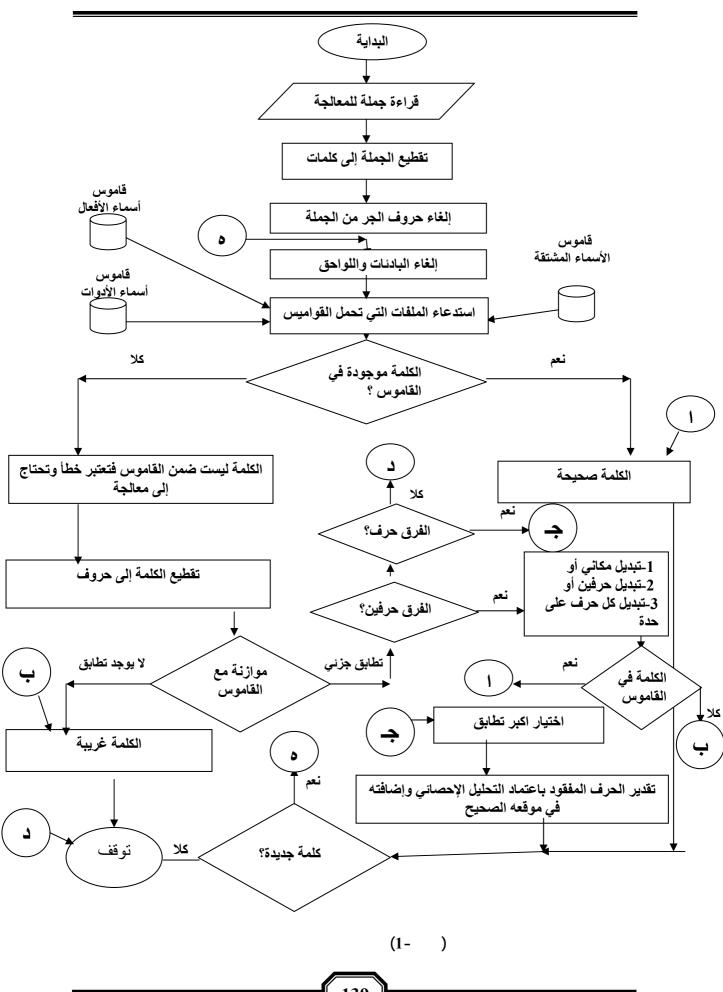
...

←

(50) (42)

.%85

. (1-)



& &

.7

OCR

.%100

(50)

.%85

(Syntax Analysis)

.(Semantic Analysis)

(1988) . **(1**

- **2)** Jurafsky, D. and Martin, J. (2000), "Speech and Language Processing", Prentice Hall.
- **3)** Tseng, Y., and Oard, D. (2001), "Document Image Retrieval Techniques for Chinese", In Symposium On Document Image Understanding Technology, pp. 151-158, Columbia, MD.
- 4) Darwish K., Hassan H., and Emam O., (2005), "Examining The Effect Of Improved Context Sensitive Morphology On Arabic Information Retrieval", Proc. Of ACL Workshop on Computational approaches to Semitic-Languages", Ann-Arbor, 2005.
- 5) Ahmad H., Sara N., and Hany H., (2008), "Language Independent Text Correction Using Finite state Automata", Proc. Of Int. Joint Conference on Natural Language Processing (IJCNLP08).
- 6) A. Amin and J.F. Mari, (1989), "Machine Recognition and correction of Printed Arabic Text", IEEE trans. On Systems, Man and Cybernities V 19, No.5, pp 1300-1306.

- 7) Sari T., and Sellami M., (2002), "MOrpho-LEXical Analysis for Correcting OCR-Generated Arabic Words (MOLEX)", Proc. Of 8th International Workshop on Frontiers in Handwriting Recognition, PP.461-466.
- 8) Walid M. and Kareem D., (2008), "Effect of OCR Error Correction On Arabic Retrieval", Information Retrieval (11), PP: 405–425.
- 9) Al-Talib Ghayda, (2006), "Fuzzy Logic Based Arabic Optical Character Recognition", Ph.D. Thesis, College of Computers and Mathematics Science, University of Mosul.
- 10) Miled H., Oliver C., et. al., (1997), "Coupling Observation / Letter For Markovian Modelisation Applied To The Recognition Of Arabic Handwriting", IEEE, pp. 580-583.
- 11) Khedher M. Z. and Abandah G. A., (2002), "Arabic character recognition using approximate stroke Sequence", 3rd Int. Conf. on Language Resources and Evaluation (LREC), Workshop.
- 12) Toufik S. and Mokhtar S., (2005), "Correcting Arabic OCR Output by Morphological Analysis of Words", International Conference on Machine Intelligence (ICMI), Tunisia.
- **13)** Wells C., Evett L., et. al., (1990), "Fast Dictionary Look-Up for Textual Word Recognition", Pattern Recognition, Vol. 23, No. 5, PP 501-508.