

The copyright © of this thesis belongs to its rightful author and/or other copyright owner. Copies can be accessed and downloaded for non-commercial or learning purposes without any charge and permission. The thesis cannot be reproduced or quoted as a whole without the permission from its rightful owner. No alteration or changes in format is allowed without permission from its rightful owner.



**HYBRID MODEL OF POST-PROCESSING TECHNIQUES FOR
ARABIC OPTICAL CHARACTER RECOGNITION**



DOCTOR OF PHILOSOPHY
UNIVERSITI UTARA MALAYSIA
2016



Awang Had Salleh
Graduate School
of Arts And Sciences

Universiti Utara Malaysia

PERAKUAN KERJA TESIS / DISERTASI
(*Certification of thesis / dissertation*)

Kami, yang bertandatangan, memperakukan bahawa
(*We, the undersigned, certify that*)

IMAD QASIM HABEEB

calon untuk Ijazah _____
(*candidate for the degree of*) PhD

telah mengemukakan tesis / disertasi yang bertajuk:
(*has presented his/her thesis / dissertation of the following title*):

**"HYBRID MODEL OF POST-PROCESSING TECHNIQUES FOR ARABIC
OPTICAL CHARACTER RECOGNITION"**

seperti yang tercatat di muka surat tajuk dan kulit tesis / disertasi.
(*as it appears on the title page and front cover of the thesis / dissertation*).

Bahawa tesis/disertasi tersebut boleh diterima dari segi bentuk serta kandungan dan meliputi bidang ilmu dengan memuaskan, sebagaimana yang ditunjukkan oleh calon dalam ujian lisan yang diadakan pada : 28 Julai 2016.

*That the said thesis/dissertation is acceptable in form and content and displays a satisfactory knowledge of the field of study as demonstrated by the candidate through an oral examination held on:
July 28, 2016.*

Pengerusi Viva:
(Chairman for VIVA)

Assoc. Prof. Dr. Haslina Mohd

Tandatangan
(Signature)

Pemeriksa Luar:
(External Examiner)

Assoc. Prof. Dr. Shahnorbanun Sahran

Tandatangan
(Signature)

Pemeriksa Dalam:
(Internal Examiner)

Assoc. Prof. Dr. Faudziah Ahmad

Tandatangan
(Signature)

Nama Penyelia/Penyelia-penyalia: Dr. Shahrul Azmi Mohd Yusof
(Name of Supervisor/Supervisors)

Tandatangan
(Signature)

Nama Penyelia/Penyelia-penyalia: Assoc. Prof. Dr. Yuhanis Yusof
(Name of Supervisor/Supervisors)

Tandatangan
(Signature)

Tarikh:
(Date) July 28, 2016

Permission to Use

In presenting this thesis in fulfillment of the requirements for a postgraduate degree from Universiti Utara Malaysia, I agree that the Universiti Library may make it freely available for inspection. I further agree that permission for the copying of this thesis in any manner, in whole or in part, for the scholarly purpose may be granted by my supervisor(s) or, in their absence, by the Dean of Awang Had Salleh Graduate School of Arts and Sciences. It is understood that any copying or publication or use of this thesis or parts thereof for financial gain shall not be allowed without my written permission. It is also understood that due recognition shall be given to me and to Universiti Utara Malaysia for any scholarly use which may be made of any material from my thesis.

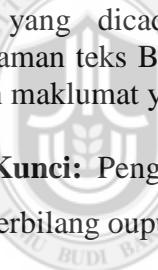
Requests for permission to copy or to make other use of materials in this thesis, in whole or in part should be addressed to:



Abstrak

Pengecaman aksara optik (OCR) digunakan untuk mengeluarkan teks yang terkandung di dalam sesuatu imej. Salah satu fasa dalam OCR ialah prapemprosesan dan ianya membentulkan kesalahan teks yang terasil dari OCR. Kaedah berbilang output dalam OCR mengandungi tiga proses iaitu: pembezaan, penajaran dan pengundian. Teknik pembezaan yang sedia ada mengalami kehilangan ciri-ciri penting kerana ia menggunakan N-versi imej sebagai input. Dalam pada itu, teknik penajaran yang terdapat dalam kajian adalah berdasarkan penghampiran manakala proses pengundian adalah tidak peka kepada konteks. Kekangan-kekangan ini mengakibatkan kadar ralat yang tinggi dalam OCR. Kajian ini telah mencadangkan tiga teknik pembezaan, penajaran dan pengundian yang ditambahbaik untuk mengatasi kekurangan yang telah dikenalpasti;. Kesemua teknik ini kemudiannya digabungkan dalam satu model hibrid yang boleh mengecam aksara optik dalam Bahasa Arab. Setiap teknik yang dicadangkan telah dibandingkan dengan tiga teknik berkaitan yang sedia ada secara berasingan. Ukuran prestasi yang digunakan adalah kadar ralat perkataan (WER), kadar ralat aksara (CER) dan kadar ralat bukan perkataan (NWER). Keputusan eksperimen menunjukkan pengurangan relatif kadar ralat dalam semua ukuran untuk teknik-teknik yang telah dinilai. Secara yang serupa, model hibrid juga telah memperolehi nilai WER, CER dan NWER yang lebih rendah iaitu sebanyak 30.35%, 52.42% dan 47.86% apabila dibandingkan dengan tiga model relevan yang sedia ada. Kajian ini menyumbang kepada domain OCR kerana model hibrid yang dicadangkan bagi teknik pasca pemprosesan boleh membantu pengecaman teks Bahasa Arab secara automatik. Oleh itu, ia akan menjurus kepada capaian maklumat yang lebih baik.

Kata Kunci: Pengecaman aksara optic Bahasa Arab, teknik pasca pemprosesan, OCR berbilang output.

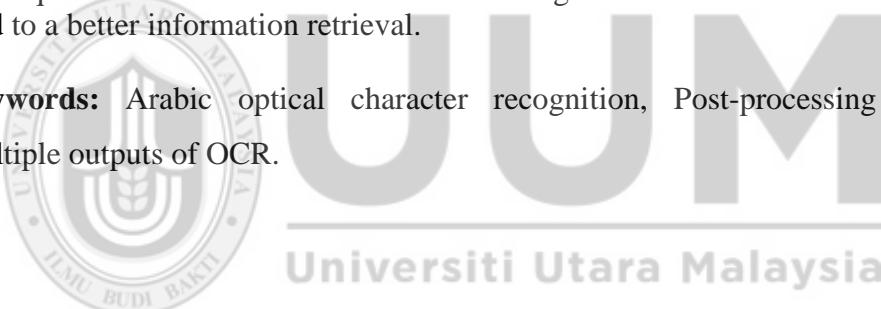


Universiti Utara Malaysia

Abstract

Optical character recognition (OCR) is used to extract text contained in an image. One of the stages in OCR is the post-processing and it corrects the errors of OCR output text. The OCR multiple outputs approach consists of three processes: differentiation, alignment, and voting. Existing differentiation techniques suffer from the loss of important features as it uses N-versions of input images. On the other hand, alignment techniques in the literatures are based on approximation while the voting process is not context-aware. These drawbacks lead to a high error rate in OCR. This research proposed three improved techniques of differentiation, alignment, and voting to overcome the identified drawbacks. These techniques were later combined into a hybrid model that can recognize the optical characters in the Arabic language. Each of the proposed technique was separately evaluated against three other relevant existing techniques. The performance measurements used in this study were Word Error Rate (WER), Character Error Rate (CER), and Non-word Error Rate (NWER). Experimental results showed a relative decrease in error rate on all measurements for the evaluated techniques. Similarly, the hybrid model also obtained lower WER, CER, and NWER by 30.35%, 52.42%, and 47.86% respectively when compared to the three relevant existing models. This study contributes to the OCR domain as the proposed hybrid model of post-processing techniques could facilitate the automatic recognition of Arabic text. Hence, it will lead to a better information retrieval.

Keywords: Arabic optical character recognition, Post-processing techniques, Multiple outputs of OCR.



Acknowledgement

Each part of this study is guided, inspired, and supported by many people. Firstly, I would like to thank all the members of my family especially my mother for their unconditional support. My goal would not be achieved without them. The most important support and guidance were from my research supervisors Dr. Shahrul Azmi Mohd Yusof and Assoc. Prof. Dr. Yuhani Binti Yusof. Thank you very much for your great help and support. It is an honor for me to do a research under your supervisions. I would like to thank all the academic and technical staff in Utara Universiti Malaysia for their help in the study process and providing all the excellent facilities. Finally, I would like to thank the Ministry of Higher Education and Scientific Research in Iraq for financial sponsorship.



Table of Contents

Permission to Use.....	i
Abstrak	ii
Abstract	iii
Acknowledgement.....	iv
Table of Contents	v
List of Tables.....	ix
List of Figures	x
Glossary of Term.....	xii
List of Abbreviations.....	xiii
CHAPTER ONE INTRODUCTION	1
1.0 Background	1
1.1 Problem Statement	8
1.2 Research Questions	11
1.3 Research Objectives	11
1.4 Significance of the Research	12
1.5 Scope of the Research	13
1.6 Organization of the Research	14
CHAPTER TWO LITERATURE REVIEW	16
2.0 Introduction	16
2.1 Arabic OCR.....	16
2.1.1 Overview of the Arabic Language	17
2.1.2 Arabic OCR Limitations	17
2.1.3 Characteristics of the Arabic Language	18
2.2 OCR Post-Processing Stage (PPS)	22
2.2.1 OCR PPS Error	22
2.2.2 Functions of OCR PPS Techniques	24
2.2.3 Categories of the OCR PPS Correction.....	25
2.3 OCR PPS Techniques	25
2.3.1 Multiple Outputs OCR (MO)	25
2.3.1.1 Differentiation Process.....	26
2.3.1.2 Alignment Process	29
2.3.1.3 Voting Process	32

2.3.2 N-grams Language Model.....	34
2.3.2.1 N-grams Language Model Functions.....	34
2.3.2.2 N-grams Language Model for Arabic	37
2.3.3 Levenshtein Distance	38
2.3.4 Rules-Based Technique.....	40
2.3.5 Noisy Channel Model	42
2.3.6 N-gram Distance	43
2.3.7 Lexicon.....	45
2.4 Comparison of OCR Post-processing Techniques	46
2.5 Hybrid Techniques of OCR PPS	48
2.6 Summary	51
CHAPTER THREE RESEARCH METHODOLOGY	52
3.0 Introduction	52
3.1 Research Phases	52
3.2 Theoretical Study	53
3.3 Design Phase	54
3.3.1 Differentiation Technique	54
3.3.2 Alignment Technique.....	58
3.3.3 Voting Technique	61
3.3.4 Hybrid Model	64
3.4 Development Phase	65
3.5 Evaluation	66
3.5.1 Data Collection.....	67
3.5.1.1 Testing Dataset.....	67
3.5.1.2 Training Dataset	68
3.5.2 Experimental Design.....	69
3.5.2.1 Differentiation Technique Evaluation.....	69
3.5.2.2 Alignment Technique Evaluation	70
3.5.2.3 Voting Technique Evaluation	72
3.5.2.4 Hybrid Model Evaluation.....	73
3.5.3 Measurments	74
3.5.4 Statistical Test	75
3.6 Summary	76

CHAPTER FOUR PROPOSED DIFFERENTIATION TECHNIQUE.....	78
4.0 Introduction	78
4.1 Differentiation Technique (EDT) Concept	78
4.2 EDT Algorithm	84
4.3 Experimental Results	86
4.3.1 Word Error Rate (WER)	86
4.3.2 Character Error Rate (CER)	89
4.3.3 Non-Word Error Rate (NWER)	91
4.3.4 Results Discussion	94
4.4 Summary	95
CHAPTER FIVE PROPOSED ALIGNMENT TECHNIQUE	96
5.0 Introduction	96
5.1 Alignment Technique (AWS) Concept	96
5.2 AWS Algorithm	100
5.3 AWS Contributions.....	102
5.4 Experimental Results	103
5.4.1 Word Error Rate (WER)	103
5.4.2 Character Error Rate (CER)	106
5.4.3 Non-Word Error Rate (NWER)	108
5.4.4 Results Discussion	111
5.5 Summary	112
CHAPTER SIX PROPOSED VOTING TECHNIQUE	113
6.0 Introduction	113
6.1 Voting Technique (VCI) Concept	113
6.2 VCI Algorithm	115
6.3 VCI Contributions	117
6.4 Experimental Results	119
6.4.1 Word Error Rate (WER)	119
6.4.2 Character Error Rate (CER)	122
6.4.3 Non-Word Error Rate (NWER)	124
6.4.4 Results Discussion	126
6.5 Summary	127
CHAPTER SEVEN PROPOSED HYBRID MODEL.....	129

7.0 Introduction	129
7.1 Interaction in the Hybrid Model (HMNL)	129
7.2 Arabic Challenges	132
7.2.1 N-gram Language Model Challenges	132
7.2.2 Diacritics	140
7.3 Experimental Results	141
7.3.1 Word Error Rate (WER)	141
7.3.2 Character Error Rate (CER)	144
7.3.3 Non-Word Error Rate (NWER)	147
7.3.4 Results Discussion	149
7.4 Summary	150
CHAPTER EIGHT CONCLUSION.....	151
8.0 Introduction	151
8.1 Achievement	151
8.2 Research Contributions	152
8.2 Research Limitations.....	154
8.3 Future Work	154
8.3 Summary	155
REFERENCES.....	157

List of Tables

Table 1.1 Some characteristics of Arabic language	3
Table 2.1 Shapes of some diacritics in Arabic	20
Table 2.2 Differentiation techniques in multiple outputs of OCR.....	26
Table 2.3 Voting techniques in multiple outputs of OCR.....	32
Table 2.4 Limitations of the OCR post-processing techniques.	47
Table 2.5 Some techniques used in the OCR post-processing stage.....	49
Table 3.1 Major variables resulted from ANOVA.....	76
Table 4.1 Experimental results of the EDT evaluation using the WER metric.	86
Table 4.2 Experimental results of the EDT evaluation using the CER metric.	89
Table 4.3 Experimental results of the EDT evaluation using the NWER metric	92
Table 5.1 Comparison between AWS technique and other existing techniques.....	102
Table 5.2 Experimental results of the AWS evaluation using the WER metric.	103
Table 5.3 Experimental results of the AWS evaluation using the CER metric.	106
Table 5.4 Experimental results of the AWS evaluation using the NWER metric.	109
Table 6.1 Voting process example.....	114
Table 6.2 Comparison between VCI technique and other existing techniques.....	118
Table 6.3 Experimental results of the VCI evaluation using the WER metric	119
Table 6.4 Experimental results of the VCI evaluation using the CER metric.	122
Table 6.5 Experimental results of the VCI evaluation using the NWER metric.	124
Table 7.1 Special tokens in the classification stage	136
Table 7.2 Example of how to store sentences in Unigram table.....	137
Table 7.3 Example of how to store sentences in Bigram table	138
Table 7.4 Example of how to store sentences in Trigram table.	138
Table 7.5 Type and size of columns of tables in N-gram language model.....	138
Table 7.6 Comparison between three Arabic corpora.....	139
Table 7.7 Experimental results of the HMNL evaluation using the WER metric.	142
Table 7.8 Experimental results of the HMNL evaluation using the CER metric.	144
Table 7.9 Experimental results of the HMNL evaluation using the NWER metric. ..	147

List of Figures

Figure 1.1. The input and output of OCR system.	1
Figure 1.2. Categories of OCR systems.	2
Figure 1.3. Stages of OCR system with output of each stage.	5
Figure 1.4. Multiple outputs of OCR.	6
Figure 1.5. Alignment process.	7
Figure 1.6. The scope of this research.	14
Figure 2.1. Connectivity in Arabic writing	19
Figure 2.2. Overlapping in Arabic writing.....	19
Figure 2.3. Diacritics in Arabic writing	21
Figure 2.4. Multiple Thresholds technique.	28
Figure 2.5. Alignment process.	30
Figure 2.6. Simple example on alignment process.	31
Figure 2.7. Levenshtein distance example.	39
Figure 2.8. Noisy channel model.	43
Figure 2.9. Bigram distance example.....	44
Figure 3.1. Research phases.	53
Figure 3.2. Flowchart of Multiple Thresholds technique.....	55
Figure 3.3. Flowchart of the proposed differentiation technique (EDT).	57
Figure 3.4. Flowchart of the existing alignment technique.....	58
Figure 3.5. Simple example of character alignment algorithm.....	58
Figure 3.6. Flowchart of the proposed alignment technique (AWS).	60
Figure 3.7. Flowchart of the existing voting technique.	62
Figure 3.8. Flowchart of the proposed voting technique (VCI).....	63
Figure 3.9. Whole evaluation process.	66
Figure 3.10. Sample image selected from the testing dataset.	68
Figure 3.11. Experiments used to evaluate the proposed differentiation technique..	69
Figure 3.12. Experiments used to evaluate the proposed alignment technique.	71
Figure 3.13. Experiments used to evaluate the proposed voting technique.	72
Figure 3.14. Experiments used to evaluate the proposed hybrid model.	73
Figure 4.1. Differentiation function and its implementation.....	79
Figure 4.2. Simple example on differentiation cycle for a primary starting pixel....	81
Figure 4.3. Simple example of proposed differentiation technique.....	83

Figure 4.4. Clustered column graph for the WER values listed in Table 4.1	87
Figure 4.5. ANOVA-test results for the WER values.....	88
Figure 4.6. Clustered column graph for the CER values listed in Table 4.2	89
Figure 4.7. ANOVA-test results for the CER values	91
Figure 4.8. Clustered column graph for the NWER values listed in Table 4.3.	92
Figure 4.9. ANOVA-test results for the NWER values	93
Figure 5.1. Loss of words' locations in MO of OCR.....	97
Figure 5.2. Extraction of words' images in the existing techniques	98
Figure 5.3. Extraction of words' images in the proposed technique.....	98
Figure 5.4. Clustered column graph for the WER values listed in Table 5.2	104
Figure 5.5. ANOVA-test results for the WER values.....	105
Figure 5.6. Clustered column graph for the CER values listed in Table 5.3	106
Figure 5.7. ANOVA-test results for the CER values	108
Figure 5.8. Clustered column graph for the NWER values listed in Table 5.4	109
Figure 5.9. ANOVA-test results for the NWER values	110
Figure 6.1. Clustered column graph for the WER values listed in Table 6.3	120
Figure 6.2. ANOVA-test results for the WER values	121
Figure 6.3. Clustered column graph for the CER values listed in Table 6.4	122
Figure 6.4. ANOVA-test results for the CER values	123
Figure 6.5. Clustered column graph for the NWER values listed in Table 6.5	125
Figure 6.6. ANOVA-test results for the NWER values	126
Figure 7.1. The interaction in the proposed hybrid model.....	130
Figure 7.2. Extract Arabic text from Wikipedia database.....	134
Figure 7.3. Database structure of N-gram language model	136
Figure 7.4. Clustered column graph for the WER values listed in Table 7.7	142
Figure 7.5. ANOVA-test results for the WER values.....	143
Figure 7.6. Clustered column graph for the CER values listed in Table 7.8	145
Figure 7.7. ANOVA-test results for the CER values	146
Figure 7.8. Clustered column graph for the NWER values listed in Table 7.9	147
Figure 7.9. ANOVA-test results for the NWER values	148

Glossary of Term

Symbol: represents the smallest meaningful unit in a writing system, such as character, number, comma, signs, etc.

Token: a sequential group of symbols not containing any spaces. It consists of a small number of symbols.

String: a sequential group of symbols. It can consist of a large number of symbols including spaces.

Word: a token exists in the specific language.

Cursive Token: a token has a group of characters joined together.

Non-word error: occurs when the word produced from the OCR process does not exist in the language resource.

Real word error: occurs when the word produced from the OCR process exists in the language resource, but it does not match with the source text.

Wrong-word: also known as an incorrect word. It refers to either non-word error or real word error.

Document Image: represents any image containing a text.

Model: a symbolic representation of concepts. It can be a schematic model or mathematical.

Lexicon: a list of words that belongs to a specific language. It does not contain any information to describe the words.

List of Abbreviations

OCR	Optical character recognition
HR	Handwriting recognition
MO	Multiple outputs of OCR
LD	Levenshtein distance
PCA	ProbCons alignment
SWA	Smith–Waterman alignment
LDB	Levenshtein distance with backtrack
WER	Word error rate
CER	Character error rate
NWER	Non-word error rate
EDT	Enhanced differentiation technique
ASW	Alignment by using words separation
VCI	Voting by using context information of sentences
CGLL	Candidates' list generation by using N-gram language model and LD
HMNL	A hybrid of MO, N-gram language model, and Levenshtein distance.
MOUMT	Multiple outputs using multiple threshold values
MOUMS	Multiple outputs using multiple scanning
MOUMO	Multiple outputs using multiple OCR systems
PPS	Post-processing stage

CHAPTER ONE

INTRODUCTION

1.0 Background

An optical character recognition, commonly referred to as OCR, is used to extract and recognize texts within images (Bassil & Alwani, 2012c). Several commercial OCR systems are currently available for various purposes, such as mail sorting systems, plate number recognition systems (Singh, Bacchuwar, & Bhasin, 2012).

Figure 1.1 shows the input and output of an OCR system.

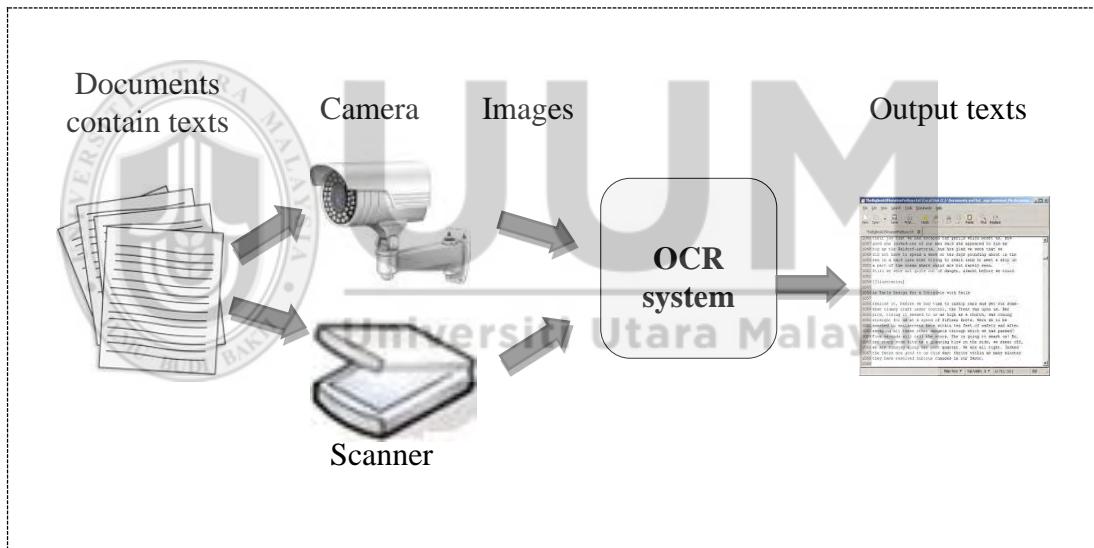


Figure 1.1. The input and output of an OCR system

There are four categories of OCR systems (El-Mahallawy, 2008). The first category is based on the type of input to these systems: offline or online. The second category depends on the mode of writing: handwritten or machine printed. The third category depends on the connectivity of a text: isolated symbols or cursive words. The last category depends on font restrictions: single font or Omni-font (Al-Badr &

The contents of
the thesis is for
internal user
only

REFERENCES

- AbdelRaouf, A., Higgins, C. A., Pridmore, T., & Khalil, M. (2010). Building a multi-modal Arabic corpus (MMAC). *International Journal on Document Analysis and Recognition (IJDAR)*, 13(4), 285-302.
- Abdulkader, A. E., & Casey, M. R. (2015). Efficient identification and correction of optical character recognition errors through learning in a multi-engine environment: Google Patents.
- Abulnaja, O. A., & Batawi, Y. A. (2012). Improving Arabic Optical Character Recognition Accuracy Using N-Version Programming Technique. *Canadian Journal on Image Processing and Computer Vision*, 3(2), 44-46.
- Ahmad, I., Mahmoud, S. A., & Fink, G. A. (2016). Open-vocabulary recognition of machine-printed Arabic text using hidden Markov models. *Pattern recognition*, 51, 97-111.
- Akhter, S., & Roberts, J. (2006). *Multi-core programming* (Vol. 33): Intel press Hillsboro.
- Akila, G., El-Menisy, M., Khaled, O., Sharaf, N., Tarhony, N., & Abdennadher, S. (2015). Kalema: Digitizing Arabic Content for Accessibility Purposes Using Crowdsourcing. In A. Gelbukh (Ed.), *Computational Linguistics and Intelligent Text Processing* (Vol. 9042, pp. 655-662): Springer International Publishing.
- Al-Badr, B., & Mahmoud, S. A. (1995). Survey and bibliography of Arabic optical text recognition. *Signal processing*, 41(1), 49-77.
- Al-Masoudi, A. F. R., & Al-Obeidi, H. S. R. (2015). Smoothing Techniques Evaluation of N-gram Language Model for Arabic OCR Post-processing. *Journal of Theoretical and Applied Information Technology*, 82(3), 432-439.
- AL-Shatnawi, A. M., AL-Salaimeh, S., AL-Zawaideh, F. H., & Omar, K. (2011). Offline arabic text recognition—an overview. *World of Computer Science and Information Technology Journal (WCSIT)*, 1(5), 184-192.
- Al-Thubaity, A. O. (2015). A 700M+ Arabic corpus: KACST Arabic corpus design and construction. *Language Resources and Evaluation*, 49(3), 721-751. doi: 10.1007/s10579-014-9284-1
- Al-Zaydi, Z. Q., & Salam, H. (2015). Multiple Outputs Techniques Evaluation for Arabic Character Recognition. *International Journal of Computer Techniques (IJCT)*, 2(5), 1-7.
- Al Azawi, M. (2015). *Statistical Language Modeling for Historical Documents using Weighted Finite-State Transducers and Long Short-Term Memory*. (PhD dissertation), Technical University of Kaiserslautern, Kaiserslautern, Germany.

- Al Azawi, M., & Breuel, T. M. (2014). *Context-dependent confusions rules for building error model using weighted finite state transducers for OCR post-processing*. Paper presented at the Proceeding of the 11th IAPR International Workshop on Document Analysis Systems (DAS) Loire Valley, France.
- Alex, B., Grover, C., Klein, E., & Tobin, R. (2012). *Digitised Historical Text: Does it have to be mediOCRe?* Paper presented at the Proceeding of the 11th Conference on Natural Language Processing (KONVENS), Vienna, Austria.
- Aljarrah, I., Al-Khaleel, O., Mhaidat, K., Alrefai, M. a., Alzu'bi, A., & Rabab'ah, M. (2012). Automated System for Arabic Optical Character Recognition with Lookup Dictionary. *Journal of Emerging Technologies in Web Intelligence*, 4(4), 362-370.
- Alkhalifa, M., & Rodríguez, H. (2009). *Automatically extending NE coverage of Arabic WordNet using Wikipedia*. Paper presented at the Proceeding of the 3rd International Conference on Arabic Language Processing (CITALA2009), Rabat, Morocco.
- Alobaedy, M. M. T. (2015). *Hybrid Ant Colony System Algorithm For Static And Dynamic Job Scheduling In Grid Computing*. (PhD thesis), Universiti Utara Malaysia, Kedah, Malaysia.
- Andoni, A., & Krauthgamer, R. (2012). The smoothed complexity of edit distance. *ACM Transactions on Algorithms (TALG)*, 8(4), 44.
- Attia, M., Rashwan, M., & Khallaaf, G. (2002). *On stochastic models, statistical disambiguation, and applications on Arabic NLP problems*. Paper presented at the Proceedings of the 3rd Conference on Software Language Engineering (CLE'2002), Cairo, Egypt.
- Attia, M., Toral, A., Tounsi, L., Monachini, M., & van Genabith, J. (2010). *An automatically built Named Entity lexicon for Arabic*. Paper presented at the Proceedings of the 7th conference on International Language Resources and Evaluation (LREC 2010) Valletta, Malta.
- Attia, M. E. (2000). *A large-scale computational processor of the Arabic morphology*. (Master thesis), Cairo University, Cairo, Egypt.
- Badawi, E.-S. M. (1996). *Understanding Arabic: essays in contemporary Arabic linguistics in honor of El-Said Badawi*: American Univ in Cairo Press.
- Bard, G. V. (2007). *Spelling-error tolerant, order-independent pass-phrases via the Damerau-Levenshtein string-edit distance metric*. Paper presented at the Proceedings of the fifth Australasian symposium on ACSW frontiers, Darlinghurst, Australia.
- Barnes, D. N. (2011). *The Text Contains its Own Lexicon: Extracting a Spelling Reference in the Presence of OCR Errors*. (Master dissertation), The Open University, Milton Keynes, United Kingdom.

- Bassil, Y., & Alwani, M. (2012a). Context-sensitive Spelling Correction Using Google Web 1T 5-Gram Information. *Computer and Information Science*, 5(3), 37-48.
- Bassil, Y., & Alwani, M. (2012b). Ocr context-sensitive error correction based on google web 1t 5-gram data set. *arXiv preprint arXiv:1204.0188*.
- Bassil, Y., & Alwani, M. (2012c). Ocr post-processing error correction algorithm using google online spelling suggestion. *Journal of Emerging Trends in Computing and Information Sciences*, 3(1), 90-99.
- Batawi, Y., & Abulnaja, O. (2012). Accuracy Evaluation of Arabic Optical Character Recognition Voting Technique: Experimental Study. *IJECS: International Journal of Electrical & Computer Sciences*, 12(1), 29-33.
- Boyell, R. L., & Ruston, H. (1963). *Hybrid techniques for real-time radar simulation*. Paper presented at the Proceedings of the November 12-14, 1963, fall joint computer conference (AFIPS '71), Las Vegas, USA.
- Cai, X. (2013). *Approximate Sequence Alignment*. (Master thesis), Louisiana State University, Louisiana, USA.
- Daðason, J. F. (2012). *Post-Correction of Icelandic OCR Text*. (Master thesis), University of Iceland, Reykjavik, Iceland.
- Damerau, F. J. (1964). A technique for computer detection and correction of spelling errors. *Communications of the ACM*, 7(3), 171-176.
- Dehkordi, Y. H. (2014). *Incorporating User Reviews as Implicit Feedback for Improving Recommender Systems*. (Master thesis), University of Victoria, Victoria, Canada.
- Do, C. B., Mahabhashyam, M. S., Brudno, M., & Batzoglou, S. (2005). ProbCons: Probabilistic consistency-based multiple sequence alignment. *Genome research*, 15(2), 330-340.
- El-Mahallawy, M. S. M. (2008). *A large scale HMM-based omni front-written OCR system for cursive scripts*. (PhD thesis), Cairo University, Cairo, Egypt.
- Farghaly, A., & Shaalan, K. (2009). Arabic natural language processing: Challenges and solutions. *ACM Transactions on Asian Language Information Processing (TALIP)*, 8(4), 14.
- Golding, A. R., & Roth, D. (1999). A winnow-based approach to context-sensitive spelling correction. *Machine learning*, 34(1), 107-130.
- Goswami, R., & Sharma, O. (2013). A Review on Character Recognition Techniques. *International Journal of Computer Applications*, 83(7), 19-23.

- Govindan, V., & Shivaprasad, A. (1990). Character recognition—a review. *Pattern recognition*, 23(7), 671-683.
- Habash, N., & Roth, R. M. (2011). *Using deep morphology to improve automatic error detection in Arabic handwriting recognition*. Paper presented at the Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Stroudsburg, USA.
- Habib, I. Q., Yusof, S. A., & Ahmad, F. B. (2014). Two Bigrams Based Language Model for Auto Correction of Arabic OCR Errors. *International Journal of Digital Content Technology and its Applications*, 8(1), 72 - 80.
- Hadj Ameur, M. S., Moulahoum, Y., & Guessoum, A. (2015). Restoration of Arabic Diacritics Using a Multilevel Statistical Model. In A. Amine, L. Bellatreche, Z. Elberrichi, J. E. Neuhold & R. Wrembel (Eds.), *Computer Science and Its Applications* (pp. 181-192). Saida, Algeria: Springer International Publishing.
- Herceg, P., Huyck, B., Johnson, C., Van Guilder, L., & Kundu, A. (2005). *Optimizing OCR accuracy for bi-tonal, noisy scans of degraded Arabic documents*. Paper presented at the Proceedings of the International Society for Optical Engineering (SPIE) on Visual Information Processing, Florida, USA.
- Howell, D. C. (2012). *Statistical methods for psychology* (8th ed.): Cengage Learning.
- Islam, A., & Inkpen, D. (2009). *Real-word spelling correction using Google Web 1T n-gram with backoff*. Paper presented at the IEEE International Conference on Natural Language Processing and Knowledge Engineering (NLP-KE2009), Dalian, China.
- Jurafsky, D., & Martin, J. H. (2009). *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition* (2nd ed.): Pearson Education India.
- Jurafsky, D., Martin, J. H., Kehler, A., Vander Linden, K., & Ward, N. (2000). *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition* (Vol. 2): MIT Press.
- Just, W. (2001). Computational complexity of multiple sequence alignment with SP-score. *Journal of computational biology*, 8(6), 615-623.
- Kai, N. (2010). *Unsupervised Post-Correction of OCR Errors*. (Diploma thesis), Leibniz University, Hannover, Germany.
- Kanoun, S., Alimi, A. M., & Lecourtier, Y. (2011). Natural language morphology integration in off-line Arabic optical text recognition. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 41(2), 579-590.

- Kenter, T., Erjavec, T., & Fišer, D. (2012). *Lexicon construction and corpus annotation of historical language with the CoBaLT editor*. Paper presented at the Proceedings of the 6th Workshop on Language Technology for Cultural Heritage, Social Sciences, and Humanities (LaTeCH 2012), Avignon, France.
- Khorsheed, M. S. (2002). Off-line Arabic character recognition—a review. *Pattern analysis & applications*, 5(1), 31-45.
- Kittler, J., Hatef, M., Duin, R. P., & Matas, J. (1998). On combining classifiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(3), 226-239.
- Knopp, J. (2010). *Classification of named entities in a large multilingual resource using the Wikipedia category system*. (Master thesis), University of Heidelberg, Heidelberg, Baden-Württemberg, Germany.
- Kolak, O., & Resnik, P. (2005). *OCR post-processing for low density languages*. Paper presented at the Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing, Vancouver, Canada.
- Kukich, K. (1992). Techniques for automatically correcting words in text. *ACM Computing Surveys (CSUR)*, 24(4), 377-439.
- Lee, Y.-S., & Chen, H.-H. (1996). Analysis of error count distributions for improving the post-processing performance of OCR. *Communication of Chinese and Oriental Languages Information Processing Society*, 6(2), 81-86.
- Lopresti, D., & Zhou, J. (1997). Using consensus sequence voting to correct OCR errors. *Computer Vision and Image Understanding*, 67(1), 39-47.
- Lund, W. B. (2014). *Ensemble Methods for Historical Machine-Printed Document Recognition*. (PhD dissertation), Brigham Young University, Utah, USA.
- Lund, W. B., Kennard, D. J., & Ringger, E. K. (2013a). *Combining multiple thresholding binarization values to improve OCR output*. Paper presented at the Proceedings of the International Society for Optical Engineering (SPIE) on Document Recognition and Retrieval XX, San Francisco, California.
- Lund, W. B., Kennard, D. J., & Ringger, E. K. (2013b). *Why multiple document image binarizations improve OCR*. Paper presented at the Proceedings of the Workshop on Historical Document Imaging and Processing (HIP 2013), Washington, USA.
- Lund, W. B., & Ringger, E. K. (2009). *Improving optical character recognition through efficient multiple system alignment*. Paper presented at the Proceedings of the 9th ACM/IEEE-CS joint conference on Digital libraries, Austin, USA.
- Lund, W. B., & Ringger, E. K. (2011, 18-21 Sept. 2011). *Error Correction with In-Domain Training Across Multiple OCR System Outputs*. Paper presented at the

Proceedings of the 11th International Conference on Document Analysis and Recognition (ICDAR 2011), Beijing, China.

Lund, W. B., Ringger, E. K., & Walker, D. D. (2014). *How well does multiple OCR error correction generalize?* Paper presented at the Proceedings of Document Recognition and Retrieval XXI (DRR 2014), San Francisco, USA.

Lund, W. B., Walker, D. D., & Ringger, E. K. (2011). *Progressive alignment and discriminative error correction for multiple OCR engines.* Paper presented at the Proceedings of the 11th International Conference on Document Analysis and Recognition (ICDAR 2011), Beijing, China.

Ma, D., & Agam, G. (2012). *Lecture video segmentation and indexing.* Paper presented at the Proceedings of the International Society for Optical Engineering (SPIE) on Document Recognition and Retrieval XIX, California, USA.

Ma, D., & Agam, G. (2013). *A super resolution framework for low resolution document image OCR.* Paper presented at the Proceedings of the International Society for Optical Engineering (SPIE) on Document Recognition and Retrieval XX, California, USA.

Magdy, W., & Darwish, K. (2008). Effect of OCR error correction on Arabic retrieval. *Information Retrieval*, 11(5), 405-425.

Mai, B. Q. Q., Huynh, T. H., & Doan, A. D. (2014). *A study about the reconstruction of remote, low resolution mobile captured text images for OCR.* Paper presented at the Proceeding of the International Conference on Advanced Technologies for Communications (ATC 2014), Saigon, Vietnam.

Muaz, A. (2011). *Urdu Optical Character Recognition System* (Master thesis), National University of Computer & Emerging Sciences, Islamabad, Pakistan.

Naseem, T. (2004). *A Hybrid Approach for Urdu Spell Checking.* (Master thesis), National University of Computer & Emerging Sciences, Islamabad, Pakistan.

Naseem, T., & Hussain, S. (2007). A novel approach for ranking spelling error corrections for Urdu. *Language Resources and Evaluation*, 41(2), 117-128. doi: 10.1007/s10579-007-9028-6

Navarro, G. (2001). A Guided tour to approximate string matching. *ACM Computing Surveys (CSUR)*, 33(1), 31-88.

Notredame, C. (2002). Recent progress in multiple sequence alignment: a survey. *Pharmacogenomics*, 3(1), 131-144.

Patel, C., Patel, A., & Patel, D. (2012). Optical character recognition by open source OCR tool tesseract: A case study. *International Journal of Computer Applications*, 55(10), 50-56.

- Pervez, M. T., Babar, M. E., Nadeem, A., Aslam, M., Awan, A. R., Aslam, N., . . . Waheed, U. (2014). Evaluating the Accuracy and Efficiency of Multiple Sequence Alignment Methods. *Evolutionary bioinformatics online*, 10, 205-217.
- Pratt, W. K. (1991). *Digital image processing*: John Wiley & Sons, Inc.
- Raaid, A. F., & Rafid, H. S. (2015). Performance Evaluation of Smoothing Techniques for Arabic Character Recognition. *International Journal of Research in Information Technology (IJRIT)*, 3(11), 22-28.
- Ramanan, M., Ramanan, A., & Charles, E. (2014). *A performance comparison and post-processing error correction technique to OCRs for printed Tamil texts*. Paper presented at the Proceeding of the 9th International Conference on Industrial and Information Systems (ICIIS) Gwalior, India.
- Rardin, R. L., & Uzsoy, R. (2001). Experimental evaluation of heuristic optimization algorithms: A tutorial. *Journal of Heuristics*, 7(3), 261-304.
- Saber, S., Ahmed, A., Elsisi, A., & Hadhoud, M. (2016). Performance Evaluation of Arabic Optical Character Recognition Engines for Noisy Inputs. In T. Gaber, A. E. Hassanien, N. El-Bendary & N. Dey (Eds.), *The 1st International Conference on Advanced Intelligent System and Informatics (AISI2015), November 28-30, 2015, Beni Suef, Egypt* (Vol. 407, pp. 449-459): Springer International Publishing.
- Sattar, S. A. (2009). *A Technique for the Design and Implementation of an OCR for Printed Nastaliue Text*. (PhD thesis), NED University of Engineering & Technology, Karachi, Pakistan.
- Shaalan, K., Samih, Y., Attia, M., Pecina, P., & van Genabith, J. (2012). Arabic Word Generation and Modelling for Spell Checking. *Language Resources and Evaluation (LREC)*, 719-725.
- Shafii, M. (2014). *Optical Character Recognition of Printed Persian/Arabic Documents*. (Doctoral dissertation), University of Windsor, Ontario, Canada.
- Shahrour, A., Khalifa, S., & Habash, N. (2015). *Improving Arabic Diacritization through Syntactic Analysis*. Paper presented at the Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, Lisbon, Portugal.
- Shannon, C., & Weaver, W. (2002). *A Mathematical Theory of Communication*: University of Illinois Press.
- Silfverberg, M., & Rueter, J. (2015). *Can Morphological Analyzers Improve the Quality of Optical Character Recognition?* Paper presented at the Proceeding of 1st International Workshop in Computational Linguistics for Uralic Languages (IWCLUL 2015), Tromsø, Norway.

- Singh, A., Bacchuwar, K., & Bhasin, A. (2012). A Survey of OCR Applications. *International Journal of Machine Learning and Computing (IJMLC)*, 2, 314-318.
- Springmann, U., Najock, D., Morgenroth, H., Schmid, H., Gotscharek, A., & Fink, F. (2014). *OCR of historical printings of Latin texts: problems, prospects, progress*. Paper presented at the Proceedings of the First International Conference on Digital Access to Textual Cultural Heritage, Madrid, Spain.
- Strohmaier, C., Ringlstetter, C., Schulz, K. U., & Mihov, S. (2003). *Lexical postcorrection of OCR-results: The web as a dynamic secondary dictionary*. Paper presented at the Proceedings of the 7th International Conference on Document Analysis and Recognition (ICDAR), Edinburgh, UK.
- Taghva, K., & Stofsky, E. (2001). OCRSpell: an interactive spelling correction system for OCR errors in text. *International Journal on Document Analysis and Recognition*, 3(3), 125-137.
- Volk, M., Furrer, L., & Sennrich, R. (2011). Strategies for reducing and correcting OCR errors *Language Technology for Cultural Heritage* (pp. 3-22): Springer press.
- Vrandečić, D., Sorg, P., & Studer, R. (2011). *Language resources extracted from Wikipedia*. Paper presented at the Proceeding of the sixth international conference on Knowledge capture (K-CAP '2011), Banff, AB, Canada.
- Vu Hoang, C. D., & Aw, A. T. (2012). *An unsupervised and data-driven approach for spell checking in Vietnamese OCR-scanned texts*. Paper presented at the Proceedings of the Workshop on Innovative Hybrid Approaches to the Processing of Textual Data, Avignon, France.
- Watson, J. C. (2007). *The phonology and morphology of Arabic*: Oxford university press.
- Yaseen, M., Attia, M., Maegaard, B., Choukri, K., Paulsson, N., Haamid, S., . . . Rashwan, M. (2006). *Building annotated written and spoken Arabic LR's in NEMLAR project*. Paper presented at the Proceeding of the 5th International Conference on Language Resources and Evaluation (LREC 2006), Genoa, Italy.
- Zribi, C. B. O., & Ahmed, M. B. (2003). *Efficient automatic correction of misspelled Arabic words based on contextual information*. Paper presented at the Proceeding of the 7th International Conference on Knowledge-Based Intelligent Information and Engineering Systems (KES 2003), Oxford, UK.