

Multilingual Artificial Text Extraction and Script Identification from Video Images



Azra Batool

Enrollment #: 01-244112-003

Supervised By

Dr. Imran Ahmed Siddiqui

Department of Computer and Software Engineering

Bahria University, Islamabad Campus

2011-2014

Multilingual Artificial Text Extraction and Script Identification from Video Images



A THESIS SUBMITTED TO THE BAHRIA UNIVERSITY IN THE PARTIAL FULFILLMENT OF
REQUIREMENTS FOR THE DEGREE OF MS SOFTWARE ENGINEERING

Azra Batool

Enrollment #: 01-244112-023

Supervised By

Dr. Imran Ahmed Siddiqui

Department of Computer and Software Engineering

Bahria University Islamabad

2011-2014

CERTIFICATE OF ORIGINALITY

I certify that the intellectual contents of the thesis

*“Multilingual Artificial Text Extraction and Script
Identification from Video Images”*

is the product of my own research work except, as cited properly and accurately in the acknowledgment and references, the material taken from any source such as research papers, research journals, books, internet, etc. solely to support, elaborate, compare and extend the earlier work, Further, this work has not been submitted previously for a degree at this or any other University.

The incorrectness if the above information, if proved at any stage, shall authorize the university to cancel my degree.

Signature: _____ Dated: January 31st, 2014 .

Name of the Research student: Azra Batool .

ACKNOWLEDGMENTS

In the name of Allah, the Most Gracious, Most Beneficent and the Most Merciful.

Alhamdulillah, all praise to Allah Almighty for the strength and blessing in completing this research work. I would like to gratefully and sincerely thank my supervisor Dr. Imran Ahmed Siddiqui for his patient guidance, understanding, encouragement and excellent supervision throughout my work. I have been lucky to have a great supervisor who cared a lot about my work and extremely honored to work in his supervision. I express my gratitude to him for his continuous and quick response to my questions and queries. I am extremely thankful to him, without his constant support and help, it would have been impossible to complete this research.

I am extremely thankful to my friend Safia Shabbir, who is always there to help me and give her best suggestions in every situation for all affairs including studies. She was always there cheering me up and stood by me through good times and bad. I would also like to thank my Uncle Mr. Shabbir Hussain and Aunt Mrs. Shabbir Hussain for their affection, love and realization of being a part of their family.

Finally, I would like to express my deepest gratitude to my beloved parents, Mrs. Sakina Raza (my Amman) and specially Mrs. Muhammad Raza (my Baba) for their love, prayers, support and encouragements throughout my life and especially during my research. I would also like to thanks to my siblings for their constant support and love.

DEDICATION

This thesis is dedicated to my adorable parents Mr. Muhammad Raza and Sakina Raza for their love, endless support and encouragement throughout my life.

Also this thesis is dedicated to the great teacher; I have ever met, Dr. Imran Ahmed Siddiqui.

ABSTRACT

With the tremendous growth in the amount of multimedia data, especially videos, has increased the need for efficient indexing and retrieval techniques. In addition to the audio-visual content itself, a power tool that be employed for indexing of videos is the caption text appearing in them. An important component of textual content based video indexing and retrieval systems is the detection and extraction of text from video frames. Most of the existing text extraction system target textual occurrences in a particular script or language. We have proposed a generic multilingual text extraction system that relies on a combination of unsupervised and supervised techniques. The unsupervised approach is based on application of image analysis techniques which exploit the contrast, alignment and geometrical properties of text and identify candidate text regions in an image. Potential text regions are then validated by an Artificial Neural Network (ANN) using a set of features computed from Gray Level Co-occurrence Matrices (GLCM). Detected text regions are then binarized to segment text from the background. The script of the extracted text is finally identified using texture based features based on Local Binary Patterns (LBP). The proposed system was evaluated on video images containing textual occurrences in five different languages including English, Urdu, Hindi, Chinese and Arabic. The promising results of the experimental evaluations validate the effectiveness of the proposed system for text extraction and script identification.

TABLE OF CONTENTS

1	INTRODUCTION	1
1.1	OVERVIEW	1
1.2	BACKGROUND	1
1.3	PROBLEM STATEMENT	2
1.4	PROPOSED METHODOLOGY	3
1.5	RESEARCH CONTRIBUTIONS	3
1.6	THESIS OUTLINE	3
2	LITERATURE REVIEW	5
2.1	OVERVIEW	5
2.2	TYPES OF TEXT IN IMAGES	5
2.2.1	<i>Scene Text</i>	5
2.2.2	<i>Caption Text</i>	6
2.3	PROPERTIES OF TEXT IN IMAGES	6
2.3.1	<i>Geometry</i>	7
2.3.1.1	Text size	7
2.3.1.2	Alignment	7
2.3.2	<i>Color and intensity</i>	7
2.3.3	<i>Motion</i>	7
2.4	STEPS FOR TEXT DETECTION	7
2.4.1	<i>Detection</i>	7
2.4.2	<i>Localization</i>	8
2.4.3	<i>Tracking</i>	8
2.4.4	<i>Extraction</i>	9
2.5	APPROACHES FOR TEXT DETECTION	9
2.5.1	<i>Un-Supervised Approaches</i>	9
2.5.1.1	Edge based methods	10
2.5.1.2	Connected component based methods	11
2.5.1.3	Texture based methods	12
2.5.1.4	Color based methods	13
2.6	SUPERVISED APPROACHES	15
2.7	SCRIPT IDENTIFICATION METHODS	17
2.8	SUMMARY	19
3	ARTIFICIAL TEXT DETECTION	21
3.1	OVERVIEW	21
3.2	CHARACTERISTICS OF TEXT IN VIDEO/IMAGES	21
3.2.1	<i>Geometrical Features</i>	21
3.2.2	<i>Edges</i>	21

3.2.3	<i>Distribution of Intensity Values</i>	22
3.3	PROPOSED METHODOLOGY	23
3.3.1	<i>Text Detection</i>	27
3.3.1.1	Image Resizing and Conversion to Grayscale.....	27
3.3.1.2	Gradient Computation.....	27
3.3.1.3	Average gradient	29
3.3.1.4	Binarization	31
3.3.1.5	Morphological Processing.....	31
3.3.1.6	Foreground Density Filter	32
3.3.1.7	Geometrical Constraints.....	34
3.3.2	<i>Validation of Text Regions</i>	37
3.3.2.1	Training	37
3.3.2.1.1	Contrast:.....	39
3.3.2.1.2	Correlation	39
3.3.2.1.3	Homogeneity	40
3.3.2.1.4	Entropy	40
3.3.2.1.5	Energy.....	40
3.3.2.2	Validation of Text regions.....	42
3.3.2.3	Text Extraction.....	42
3.4	SCRIPT IDENTIFICATION	45
3.4.1	<i>Local Binary Patterns</i>	46
3.4.2	<i>Training and Classification</i>	48
3.5	SUMMARY	48
4	RESULTS AND EXPERIMENTS	50
4.1	OVERVIEW	50
4.2	DATA SET.....	50
4.3	EVALUATION METRICS	51
4.4	GROUND TRUTH LABELING.....	52
4.5	TEXT DETECTION RESULTS	54
4.6	SCRIPTS RECOGNITION RESULTS.....	56
4.7	SENSITIVITY TO SYSTEM PARAMETERS	56
4.8	SUMMARY	58
5.	CONCLUSION AND PERSPECTIVES	59
5.1.	CONCLUSION.....	59
5.2.	PERSPECTIVES	59
	BIBLIOGRAPHY	61

List of Figures

Figure 2.1: Examples of scene text appearing in videos	6
Figure 2.2: Examples of artificial text in videos	6
Figure 2.3: Examples of text detection	8
Figure 2.4: Examples of text localization (a) Input image frame (b) Text localization	8
Figure 2.5: (a) Text Localization (b) Extracted text	9
Figure 3.1: Image blocks and their intensity histograms, (a): Text on homogenous background (b): Text on complex background (c): Non-text region	23
Figure 3.2: Overview of video indexing and retrieval system	24
Figure 3.3: Overview of proposed methodology	25
Figure 3.4: Overview of steps involved in text detection and validation	26
Figure 3.5: Samples of textual content in five different languages (a) Urdu (b) Arabic (c) Chinese (d) Hindi (e) English	28
Figure 3.6: Sobel operator for detection of vertical edges	28
Figure 3.7: (a, c) Gray scale images with Chinese and Urdu text respectively, (b, d) Gradient image showing vertical edges	29
Figure 3.8: Average gradients of the images in Figure 3.7	30
Figure 3.9: Application of horizontal dilation to merge components	32
Figure 3.10: Images after morphological processing (left column) and after application of density filter (right column)	33
Figure 3.11: Localized text regions in two images	36
Figure 3.12: Images after application of geometrical constraints	36
Figure 3.13: Sample training examples (a): Non-text blocks (b) Text blocks	38
Figure 3.14: (a) An image with 4 gray levels (b) GLCM using (1, 0) displacement vector	39
Figure 3.15: A simple feed forward neural network	41
Figure 3.16: Examples of text segmented from the background	45
Figure 3.17: LBP Computation (a): Image values (b): Binary codes assignment (c): Weights of neighboring pixels (d): Conversion to decimal	46
Figure 3.18: Examples of the ELBP operator [55]. The circular (8, 1), (16, 2), and (24, 3) neighborhoods.	47
Figure 3.19: Input Text blocks	48
Figure 4.1: Sample images and corresponding ground truth images	54

Figure 4.2: (a): Detected text region (b): Ground truth text region 55

Figure 4.3: Detection performance as a function of foreground density filter threshold..... 57

Figure 4.4: Script recognition rates as a function of different neighborhoods of LBP 58

List of Tables

Table 3.1: Values of threshold on geometrical constraints	36
Table 4.1: Distribution of dataset.....	50
Table 4.2: Precision and recall of text detection (unsupervised)	54
Table 4.3. Precision and recall after text validation.....	56
Table 4.4: Script Recognition – Confusion matrix	56