

Removal of Graphics from Text-document and Segmentation of Gujarati Documents Using Connected Component Theory

Himanshukumar U. Prajapati¹, Prof. S. Ramamohan², Ms Sonam Chauhan³

¹(Faculty Of Tech & Engg., MSU, Vadodara)

²(Chief coordinator, MCIT-OCR Project, MSU, Vadodara)

³(Electrical Engg. Dept., Faculty Of Tech & Engg., MSU, Vadodara)

ABSTRACT

Very little work is found in the literature for recognition of Indian language scripts. We have work on graphics separation & text segmentation for Gujarati language documents.

Text/graphics separation aims at segmenting the document into two layers: a layer assumed to contain text and a layer containing graphical objects. In this paper, we propose method of choosing right threshold to extract text from graphics. We consider connected component characteristics to differentiate between text and graphics. It is based on the observation that connected components of graphics are comparatively larger than text. Experimental results show that while removing graphics, no damage on connected components of text in the image. Segmentation is one of the most difficult tasks in digital image processing. In which each character of the document has to be identified individually.

Keywords - connected components, graphics separation, minimum bounding box (MBB), and segmentation.

INTRODUCTION

OCR software for Gujarati Language is almost in final stage of development after the continuous hard work for many years but there are several challenges which have to be solved.

Challenges:- Gujarati OCR has been currently in development but faces many challenges that reduce its efficiency. In some kind of books, there are characters having irregular thickness. Hence, problem occurs due to Broken and merged characters in document images pose serious challenges for recognition accuracy. If these problems can be reduced by any means, then overall efficiency can be increase by up to 10-15% a now

-There are graphics related issues; Which reduces this system's efficiency.

-Recognition of handwritten document is difficult one, because different persons have their own writing style. New systems have been developed for that which is known as ICRS.

-There are problems of fusion of different languages includes in single document Like: English & Gujarati or Sanskrit & Gujarati

The document image generated as input to OCR may be subjected to various operations to remove the artifacts and enhance the image quality so that subsequent processes give expected results. Thresholding (or binarization) & graphics removal are complex challenges during pre-processing stage. Documents containing graphics is one of the challenges to develop OCR (Optical character recognition system). It reduces OCR efficiency rapidly. Graphics has not uniform shape or design, it always varies. So in our database we have no data about any particular graphics. That is why it is necessary to remove graphics for better result of OCR, output of which is editable text.

II. CHARACTERISTICS OF GUJARATI SCRIPT

As other Indian languages the character set of Gujarati comprises of 36 consonants, 12 vowel and 6 signs, 12 dependent vowel signs, 10 digits. The consonants can be combined with the vowels and can form compound characters as shown in Fig 1.

અ	આ	ઇ	ઈ	ઉ	ઊ	એ	ઐ	ઓ	ઔ	અં	અઃ
a	ā	i	ī	u	ū	e	ai	o	au	am	ah
[ə]	[ā]	[i]	[ī]	[u]	[ū]	[e, e]	[əy]	[o, ɔ]	[əu]	[əŋ]	[əh]
મ	મા	મિ	મી	મુ	મૂ	મે	મા	મો	મા	મં	મઃ
ma	mā	mi	mī	mu	mū	me	mai	mo	mau	maṅ	māh

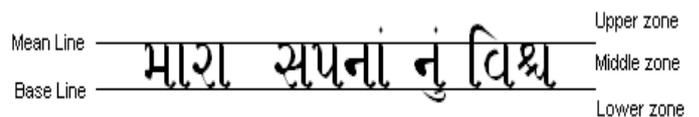
A word may be formed by combining the basic character(s), which may be combined with vowel(s). Sometimes the conjunct consonants may also form part of the word. Collection of words will form a line and collection of lines will form text.

Fig. 2 shows a line of Gujarati text, which can be considered as being logically divided into three horizontal parallel lines [3]. The base line, the imaginary line separating the middle and lower zone on which consonants and independent vowels are written (Middle zone).

The line below the base line, used for writing dependent (lower) vowels (Lower Zone).

The line above the mean line, used for writing dependent (upper) vowels (Upper zone).

A typical example of zoning is shown in Fig



ACKNOWLEDGEMENT

This work has been carried out at MCIT-OCR centre at Linguistic Dept. Arts Faculty, MSU, Vadodara. I am so thankful to Ms Soram Kotia & Prof. A. I. Trivedi at Faculty of Tech. & Engg. for their continuous guidance and precious help during work

REFERENCES

Journal Papers:

- 1) Ruini Cao, Chew Lim Tan School of Computing, "Text/Graphics Separation in Maps" National University of Singapore.3 Science Drive 2, Singapore 1175 (Text graphics separation) ,2001
- 2) Karl Tombre, Salvatore Tabbone, Loïc Pélissier, Bart Lamiroy, and Philippe Dosch, LORIA "Text/Graphics Separation Revisited" B.P. 239, 54506 Vanduvrelès-Nancy France (Text graphics separation), 2002
- 3) S.Rama Mohan, Jignesh Dholakia, Atul Negi. "Zone identification in the printed Gujarati Text" Processing of the 2005 Eight International Conference on Document Analysis & Recognition(ICDAR'05)

Books:

- 4) Gonzalez, R.C. and Woods, R.E,"Digital Image Processing", Second Edition, Pearson Education, Singapore

Theses:

- 5) Mr. Jignesh Dholakia, Mathematical Techniques for Gujarati Document Analysis and Character Recognition,2010

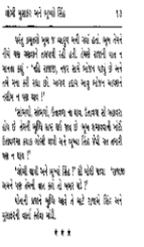
Websites:

- 6) http://en.wikipedia.org/wiki/Optical_character_recognition (Basics of OCR)
- 7) nstats.un.org/unsd/demographic/meetings/wshops/new.../docs/05a.ppt(Basics of OCR)

Text-Graphics Separation Result

		
<p>Input</p>	<p>Removed Graphics</p>	<p>Output</p>

Text-Graphics Separation Result

		
<p>Input</p>	<p>Removed Graphics</p>	<p>Output</p>