

35	ی	ی ی ی	یٹوھچ ے ی	čhōt.ī yē
35a	ئ	ئ ئ ئ	مزم	hamzah
35b	ے	ے	ے ی ٹ ب	bar.ī yē

Urdu is written in Arabic script. Arabic script has many traditional writing styles, including Naskh (mostly used for Arabic language), Taleeq, Kufi, Divani, Sulus, Rika, etc. Naskh and Taleeq styles of writing were combined into the very spatially concise Nastaleeq writing style. Nastaleeq writing system for Urdu is character based, bidirectional (mainly R to L), diagonal, non-monotonic, cursive, context sensitive writing system with a significant number of marks (dots and other diacritics). This makes Nastaleeq one of the most complex writing styles and challenging to develop an OCR for it. Nastaleeq is a complex cursive style of writing Arabic script based languages e.g. Urdu and Persian. Each letter has precise writing rules, relative to the width of the flat nib of the pen, called *qat*. The measurement of some letters in terms of *qat* is given in Fig. 2.



As Nastaleeq is a writing style for Arabic script, it inherits its bidirectional nature, where the characters are written from R to L but numbers are written from L to R. Hand written Nastaleeq has been developed as art in the Muslim world where it replaced all other forms of art like painting etc.

The main problem arising in segmentation is the possibility of overlapping of characters in a word or sub-word which occurs quite often especially in cursive languages. Eliminating the possibility of overlapping by stretching the words horizontally to make space between two connected characters is shown in Fig. 3. Text lines, words boundary identification and overlapping and other preprocessing techniques have been adopted from [1] [3], [4], [5] and [6].

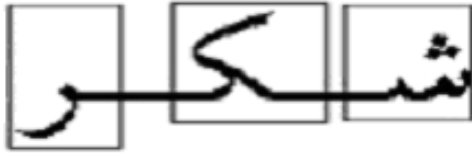


Fig. 3 Horizontally Stretched Word

B. Segmentation

The segmentation phase is based on the level of complexity offered by a character during scanning. The characters are grouped into three levels of complexity, simple, semi complex and complex as shown in Table II.

TABLE II
COMPLEXITY WISE GROUPING OF CHARACTERS

Char	Forms	Complexity Level
ف	ل اکشا	Simple
ا	ا	
ب	ب ب ب	
ٹ	ٹ ٹ ٹ	
ش	ش ش ش	Semi Complex
ء		
ج	ج ج ج	
د	د	
ڈ	ڈ	Semi Complex
ک	ک ک ک	
گ	گ گ گ	
ل	ل ل ل	
ی	ی ی ی	Semi Complex
ے	ے	
ن	ن ن ن	
ض	ض ض ض	
ظ	ظ ظ ظ	Complex
غ	غ غ غ	
ف	ف ف ف	
ق	ق ق ق	
م	م م م	Complex
و	و	

~	~ ~ ~	Complex
ھ	ھ ھ ھ	Complex

The complexity of a character is measured by analyzing the topological features, the number of holes, the width, height of holes and the direction of these holes but the decisive part is played by the lines which are encountered by the scanner

during the scanning process. The character dō-čāšmī hē (ھ) is made of two holes by three lines (connecting each other) therefore is considered a complex character similarly a single

hole or closed / loop character like mīm (م) is also considered as a complex character. A character with semi opened shape or with two line from one side like bar.ī Hē (ح) is considered semi complex shape, all the remaining characters are considered simple ones. Deciding a complex or semi complex shape it has taken care that the lines should be connected at some point and the distance of any two lines should not exceed a specified limit keeping in view the size of the fonts under consideration like the distance between the two lines of bar.ī Hē (ح) .

To avoid complex calculation and improve the efficiency the scanning is carried out both vertically and horizontally. During which an isolated word is scanned vertically from right to left, double and triple lines characters are looked firstly from the upper side of the image, if the character is closed with two or more lines such that it makes a hole it is further scanned horizontally from top to bottom to define the hole as a vertical or horizontal hole and calculate the distance between the lines of holes. If the character is closed from three sides but open or semi open from one side these characters are semi complex characters. All other characters are considered simple one.

The character is scanned and the level of complexity is stored during the scanning as complex, semi complex or simple, when the level is getting changed (the loop is starting / ending) the change is verified from the right side scan if it also confirms the change the beginning /end for complexity is marked, reaching on the other side of the character the same process is performed, now the character is marked on two ends as a single / isolated character. It is extracted from the word and the search for a new character continues.

C. Character Recognition

Character Recognition has been performed using Neural Networks; the technique used here is described in detail in our paper [6]. All the procedure is the same except that this time the training was done using different forms of a character.

V. CORPORA

For experiments, we collected a corpora consisting of two sets of images (and associated transcriptions): computer generated, i.e. synthetic, images and real-world images consisting of scans of commonly available hardcopy Urdu documents which do not contain any other language.

VI. ASSUMPTIONS MADE

During the whole process of segmentation and character recognition the input script is assumed to be diacritic (Erabs) free. The font's size has been kept fixed or the image has been resized to make the fonts suitable for segmentation and recognition.

VII. RESULTS

Old and newly written scripts were used to evaluate results on good and bad quality paper which produced results on the average as 93.4%, which can further be improved using a lexicon and focusing more onto frequently used characters in the script.

VIII. FUTURE DIRECTIONS

This Urdu character recognition system is developed on diacritic (Erabs) free Urdu text. Further research is needed to develop a system that recognize diacritic of Urdu, Arabic and other languages having the same properties, with an integrated lexicon to further improve the results.

IX. CONCLUSION

This paper describes a system for OCR of printed Urdu script. The recognition accuracy of our prototype is promising, but more work is needed. Our character segmentation method should include handling a larger variety of characters including roman script that occurs often in images obtained from Urdu documents. We also need to recognize characters having diacritic to make it a complete OCR system. In general, the system needs to be tested and fine-tuned on a wider variety of images containing characters in diverse fonts and size.

REFERENCES

- [1] U. Pal and Anirban Sarkar, "Recognition of Printed Urdu Script", "Proceedings of the Seventh International Conference on Document Analysis and Recognition (ICDAR 2003)".
- [2] Raymond G. Gordon, "Ethnologue: Languages of the World Fifteenth Edition" SIL International, 2005.
- [3] Khalid Saeed, "New Approaches for Cursive Languages Recognition: Machine and Hand Written Script and Texts".
- [4] K. Saeed, Three-Agent System for Cursive Script Recognition, " Proc. CVPRIP '2000 Computer Vision, Pattern Recognition and Image Processing-5th Joint Conf. on Information Sciences, JCIS'200, Vol.2, PP.244-247, Feb 27-March 3, N.Jersy 2000.
- [5] K. Saeed, R Niedzielski, "Experiments on Thinning of Cursive-Style Alphabets, "Inter Conf. on information Technologies ITESB '99, June 24-25, Minsk 1999.
- [6] Inam shamsheer, Zaheer Ahmad, Jehanzeb Khan Orakzai, Awais Adnan, "OCR For Printed Urdu Script Using Feed Forward Neural Network," MLPR 2007 :International Conference on Machine Learning and Pattern Recognition", 2007.