



RESEARCH ARTICLE

AUTOMATED DATA CAPTURING AND RECOGNITION TECHNOLOGY TO MINIMIZE HUMAN INTERVENTION IN UNIVERSITY EXAMINATION SYSTEM

*Mohini and Dr. Amar Jeet Singh

Department of Computer Science, Himachal Pradesh University, Shimla, India

Received 29th September, 2012; Received in Revised form; 28th October, 2012; Accepted 11th November, 2012; Published online 17th December, 2012

ABSTRACT

Manual data entry from hand written data is very time consuming, tedious and prone to many errors more so where there is bulk amount of data is involved. This is the area where there is a requirement of automated data capturing and recognition technology. Using a recognition engine to convert text or handwriting from the printed page into computer readable characters saves up to 90 percent of the time it would take to enter the information manually. To fetch the data with minimal human intervention there are three types of recognition engines for this purpose and these are Optical Character Recognition (OCR), Intelligent Character Recognition (ICR) and Optical Mark Reader (OMR). Generally, OCR is best for machine print, or type, ICR is better for converting handwriting data and OMR is best suitable to detect the absence or presence of a mark, but not the shape of the mark. Automated data capturing is rapidly becoming an integral and necessary component in any organization. This not only saves the cost but also increases the speed and accuracy over manually entered data. This paper report on the study of use of these technology in various processes of examination system especially in fetching awards from OMR technology to get almost hundred percent accuracy.

Key words: OCR, OMR, ICR, AutoRec.

INTRODUCTION

Due to technological advancement human has been successful in getting handwritten data converted into machine readable form without wasting much time and cost. Computer can understand only alphanumeric characters as American Standard Code for Information Interchange (ASC II) typed on a keyboard where each character or letter has a unique code. However, computer itself does not recognize characters and numbers from scanned image. This has been overcome by applications of Artificial Intelligence to match images of characters available on scanned documents and convert them into their ASC II equivalents to get readable text. An Intelligent Character Recognition (ICR) System for handwritten forms includes functional components for form registration, character image extraction and character image classification. An effort is needed so that written words get separated into individual characters. Using ICR technology, computer system becomes artificially intelligent to recognise characters based on their shapes. The matched characters are directly converted into machine readable form. For unrecognised or doubtful characters, human intervention is needed to convert them to correct characters. So, ICR technology is seemingly a good machine facility to human operators to minimise their data entry time, decrease human drudgery and increase overall productivity. But in tasks where 100% accuracy is desired such as awards in examination system, this technology does not give desirable results. Himachal Pradesh University, Shimla, India is using ICR technology to fetch students awards. So here, the demand is of

that technological which can give 100% accuracy. In this paper we are discussing the practical problems in data fetching from handwritten forms through ICR. Before discussing the recognition results, it is instructive to review various component of ICR technology.

Processing of ICR Forms

Examination system of any university is most important component in University Administration. Manual Examination system is very tedious job and due to human intervention it sometimes leads to errors which tarnish the image of university. Himachal Pradesh University has initiated a process to use the ICR technology, an Automated Data Capturing/Recognizing System to minimize human intervention and for fetching data through ICR, AutoRec software is being used. An *AutoRec* is an intelligent data capture and perfection software that facilitates the conversion of hard copy data into computer readable form /electronic data with minimum human interventions. The various components of "AutoRec" Flow Process are shown in Figure 1. ISCAN1 and ISCAN2 are two scanner stations where we can scan ICR forms. There are different processes to fetch the data from ICR forms, these are briefly discussed below:

Designing of Template: The first step in fetching data from ICR awards is to design a template for the document whose data has to be recognized. This is necessary to design the form according to position of the fields. Type of fields, position of data fields, length of field, type of recognition and other properties associated with any data item which we want to

recognize. Frame is designed and registration points are set and template is registered.

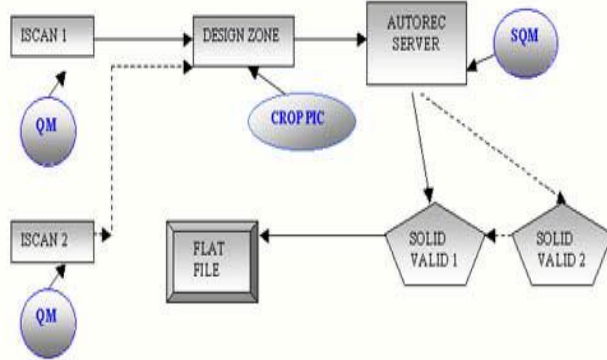


Fig. 1: AutoRec Flow Process (Source : <http://www.fil-flan.com>)

Image Capture & Scan: Different batches are created and forms are scanned. The forms which are scanned are images and are placed in different batches. We can scan more than 100 forms at a time from a scanner according to the capacity of the scanner.

Image Processing: In this process scanned image is processed. Here different processes are involved viz. image rotation, deskewness, despeckling, border removing, edge noise or entire background, locating data field and their identification and registration of the form.

Batch Processing: The batches are processed according to designed template.

Recognition: After processing the batches next step is to recognize them. The recognition starts processing one form at a time from the batch folder and marks the data fields in different colors depending on the recognition confidence. It is capable of recognizing free hand writing, hand written numbers and alphabets, check boxes, bar codes, OCR characters and magnetic ink characters. The recognized images are then passed to the data validation.

Data Validation: In this process, data is distributed for validation to different users. Operators are deployed for data validation. In data validation, there is a split screen, the actual image and the recognized data are shown separately. The user can see the image with the naked eye, check and correct any unrecognized data elements. After doing the required validation data can be exported, the data can be exported to various data formats e.g. flat ASCII file, PDF, HTML or ODBC compliant data bases.

Exported Files: The exported files are ready for data processing

Related work: Presently university is using ICR technology to fetch awards from hand written data to machine readable form. Deodhare, Dipti, N. N. R. Ranga and Suri has suggested two most significant technologies, Intelligent Character Recognition (ICR) for hand written documents and Optical Mark Recognition (OCR) for data capture from printed documents. An Intelligent Character Recognition (ICR) System for handwritten forms includes functional components for form registration, character image extraction and character image classification. Rosen, Richard, Vinod and Patrick has

suggested that ICR technology along with FAX system is a viable cost-effective option for processing CES reports presently being received via paper FAX. Overall, the cost of the ICR system would be recovered within a year. Fax ICR technology is sufficiently accurate to provide high quality data with much less operator time than is presently used for full key entry of the data. ICR mode of data fetching is not 100 percent accurate as it is not able to recognize poor handwriting. Sometime evaluators marks stray marks on the award list and rollnos are also not properly written. Viking has stated that ICR engines can achieve very high recognition rates when the documents are properly designed, printed and controlled. Nonetheless, a 1% error rate on a printed page with 3,000 characters means there is an average of 30 errors on each page. This would be unacceptable for a typist, but may be adequate for information that will only be read occasionally. Kevin has found three practical methods of getting data from a piece of paper into a computer. These three methods are to key the data, to use Optical Mark Reading (OMR) or to use Imaging combined in some way with Optical Character Recognition (OCR). These are reliable and validated transfer of completed form data to a computer system for further processing and has suggested that OMR is a fast, accurate and cheap method for collecting census data. According to Bergeron OMR systems approach one hundred percent accuracy and only take .005 seconds on average to recognize marks.

ICR Form for fetching data from Awards List for university examination : As discussed in related work, ICR is a good technology that can be used to minimize human intervention to get accuracy. But recognition of hand-printed forms can be challenging. Form design is vital to ICR accuracy. ICR reads images of hand-printed characters from paper and converts them into machine-readable characters. The ICR award list which is currently being used in HP University is shown in Fig. 2.

Fig. 2 : ICR Award List

The award list shown above is designed in red colour. To capture the data, red colour is dropped and data written in the boxes is fetched. Accuracy of data depends upon the handwriting. There are other practical problem which the evaluators make while filling the awards in the ICR forms. Although instructions have been provided but still evaluators do make mistakes. Some of the problems are shown below:

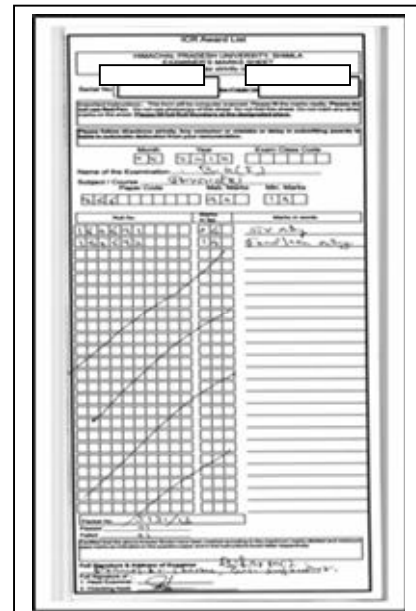
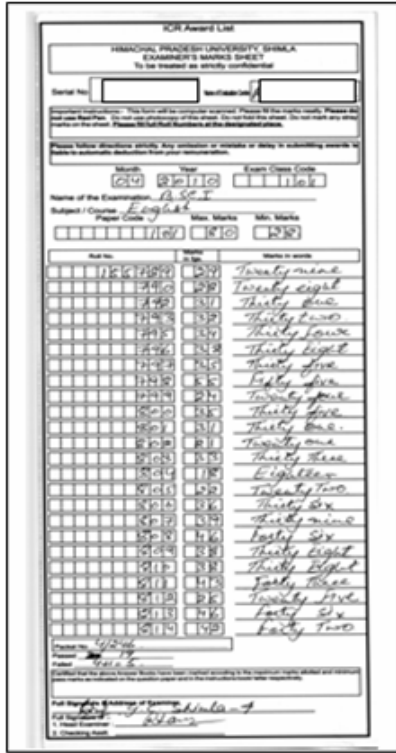


Fig. 5 : Stray marks on ICR marks list

Processing and recognition of handwritten documents is relatively difficult. The cursive scripts, variability and separation of characters, document skew, accommodating variability of strokes, learning the human characteristics are very difficult to model. Most of the handwritten document understanding systems are employed in specific narrow domain. HP University is adopting ICR technology in fetching awards of BSc. And BCom classes only. As stated above, ICR technology is not hundred percent accurate. It sometimes misinterpret the data and even it does not have doubt its recognition. The Figure below explain this fact. In the figure ICR recognize the Figure as 5594 whereas the actual figure is 5534 as shown in the scanned image.



Fig 3-4 : Complete roll nos are not written, only last three/ two digits of roll nos are written. ICR is not intelligent enough to recognize candidate's complete rollno.

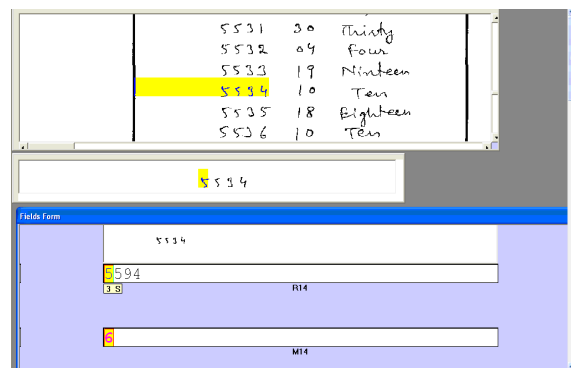


Fig 6 : Misrecognition of figure 5534 as 5594

OMR for accuracy in Awards

ICR is less accurate than OMR and it also requires some editing and verification. But examination system demands 100 percent accuracy in awards fetching. For awards fetching we must use that technology which gives more accuracy than ICR. OMR is the fastest and most accurate of the data collection technologies. It is also relatively user-friendly. OMR technology detects the existence of a mark, not its shape. According to Webster's Online Dictionary, Optical

