

Anatomy of a Form Reader

Stephen W. Lam, Ladan Javanbakht and Sargur N. Srihari
Center of Excellence for Document Analysis and Recognition
State University of New York at Buffalo
Buffalo, New York 14260

Abstract

Forms are used extensively in today's offices. Task of an automated form reader is to locate data filled on a form and to encode the content into appropriate symbolic descriptions. The challenges in form reading are due to high volume and large variety. This paper presents a robust form reader with high adaptability and trainability. The form reader consists of two modules: field registration and data recognition module. The field registration module acquires knowledge about the forms of interest and the data recognition module recognizes text data on filled forms using the acquired knowledge. Capability of the reader increases progressively through supervised learning. The form reader has been training to read a large variety of forms with machine-printed data. The adaptability and trainability of the system have been demonstrated through the experiments.

1 Introduction

A typical automated forms processing system is composed of two major modules: *form registration* and *data processing*. The form registration module detects unique topological features and symbolic codes of a form that can be used for form discrimination. The module also defines areas where data will be found. The data processing module determines the identity of a form by matching its topological and symbolic information against those in a database. Data will be extracted from the predefined areas for recognition [?]. Many commercially available form readers have very limited application since they can only read forms with specific designs, such as printed with dropout ink and fixed identity location. In many cases, forms have to be redesigned for the readers. Data areas on a form are required to define by tedious manual operation. This becomes very expensive and time consuming for large application and frequent form addition. Moreover, valuable contextual information are seldom used intelligently to improve OCR accuracy [3]. Based on these observations, a general purpose forms reader which addresses all these issues is being developed at CEDAR. Novelties of the approaches taken are able to attain high degrees of success in *adaptability*, *learnability* and *accuracy*.

Experiments are being conducted on machine-printed IRS tax forms, medical insurance claim forms and office forms. Hand-printed data will be included in future tests. No constraint is imposed on form de-

sign and a semi-automatic field registration scheme is employed to eliminate manual creation of data areas. Several word recognition strategies and natural language processing techniques are developed and they are used selectively according to data contextual constraints. A processing rate of 15 forms per minute is achievable with special hardware support.

2 System Overview

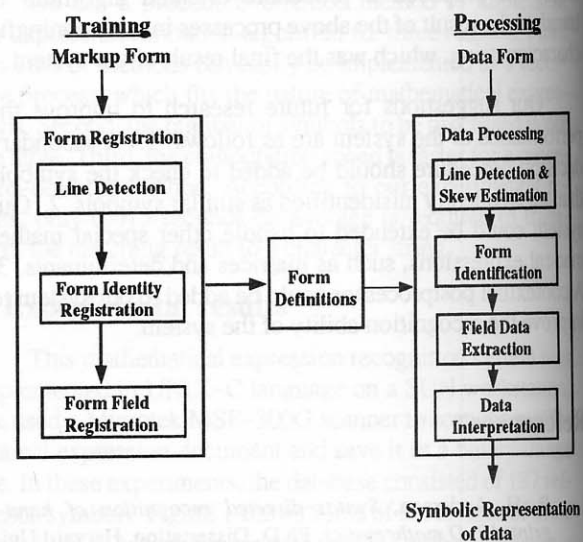


Figure 1: System components of the CEDAR form reader.

The major components of the CEDAR form readers and the overall forms processing stages are depicted in Figure 1. A blank form with markup codes typed in the data areas is presented to the form registration module for automatic form identity and field location registration. A markup code, such as "@A0B", is a syntactic string which indicates location of a data field, field boundary type and data constraint. Each character position in this syntactic string represents a specific type of information and the character is the attribute of the information. Extracted topological and symbolic information together with the field lo-

cations and contextual constraints are stored in the form definition database. The data processing module first extracts topological and symbolic information from the data form and finds the closest match in the database. If a successful match occurs, then the pre-defined field locations and constraints will be used to guide data extraction and recognition.

The line detector and form identifier are used in both modules. The line detector is also capable to estimate form skew angle which critical to data extraction and recognition.

3 Line Detection and Skew Estimation

The first stage in both form registration and data processing modules is line detection and skew estimation. Lines provides important topological features for form identification and skew information is needed to adjust field locations in data extraction. Hough transform is a popular approach to detect skew angle of a document [4]. However, this method is time-consuming to process a large image in software. Lines are detected from the contours of form components. Consecutive chaincodes on a contour that have similar direction form a line and the chaincode direction represents the line orientation. The form skew angle can easily be estimated by using majority vote on line orientation. Figure 2 shows line detection using object contours. For a letter-sized form, the current software version takes about 11 seconds on a SPARC10 to generate contours and only 0.05 second to detect lines and estimate skew angle. However, contour generation can be achieved in less than 1 second by using a contour tracing hardware board. This hardware board is being used in our real-time postal address interpretation project.

Line-touching characters can easily be detected by finding sudden change in contour direction which is not orthogonal to the horizontal line. Portion of the contour that belongs a character will be marked. Characters will not be detached from the lines until the data interpretation stage. Contextual constraint and non line-touching data of a field can provide clues in separating the characters.

4 Form Identity Registration and Recognition

Form identity registration is performed at the form registration stage. It has to be reliable under various form variations. A incorrect identification means that wrong field locations and data constraints will be used in recognizing the data. This identity classifier first captures topological features that are unique to a particular form. The features should have acceptable consistency to occur on different variations of the same form, *i.e.* they must be invariant to form rotation, shifting and translation, and have high tolerance to noise. Line intersections such as those in Figure 3 are currently used as features for classification. The number of features and their relative spatial positions should be the same on different copies of the same form. A form is spatially partitioned into 9 equal sized regions (See Figure 4). In a region, the number

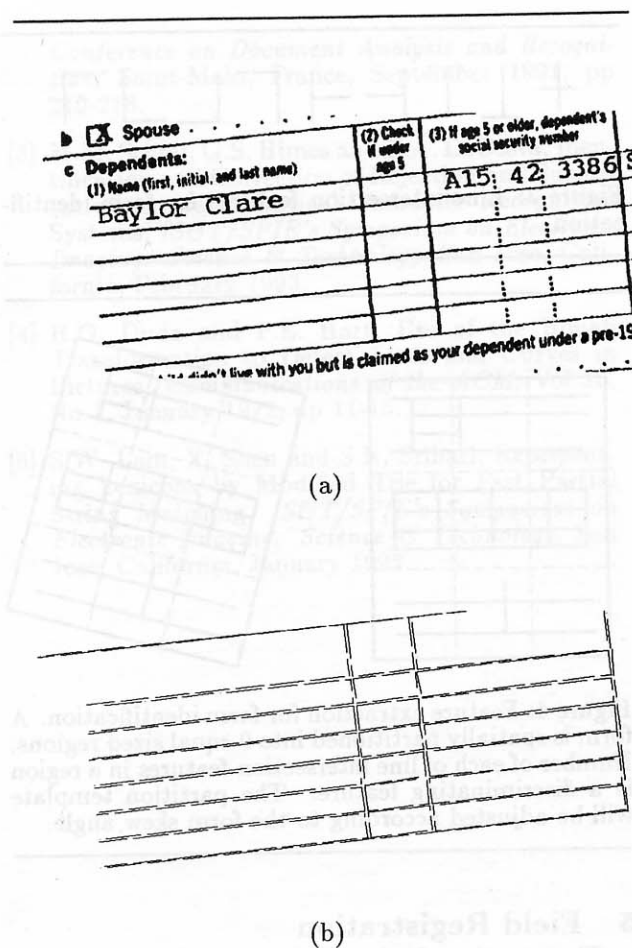


Figure 2: Stages in line detection and skew estimation. (a) Portion of a tax form. (b) Lines derived from the large contours by chaincodes tracing. 5 degrees skew is estimated based on the line orientation.

of each of the line intersections is considered as a feature. Therefore, a multi-valued feature vector which contains 81 line intersection features is formed.

Such a feature vector represents a point in the 81-dimensional feature space and the task of a classifier is to find the closest point that represents a markup form. An Adaptive Resonance Theory 1 (ART1) neural network is being used as the classifier. This classifier has the capability to self-organize a new form category when none of the feature vectors in the database has a close match with the processing form.

Nevertheless, topological features alone are not sufficient to distinguish similar forms, *e.g.* same tax form of different years. Symbolic information such as form name is necessary to make fine discrimination. Therefore, the form name will also be located at the identity registration stage and may be used at the recognition stage when topological features cannot make high confidence decision.

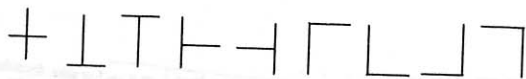


Figure 3: Line intersection features for form identification.

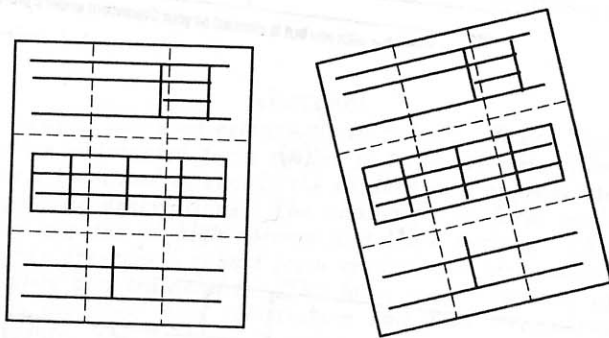


Figure 4: Feature extraction for form identification. A form is spatially partitioned into 9 equal sized regions. Number of each of line intersection features in a region is a discriminating feature. The partition template will be adjusted according to the form skew angle.

5 Field Registration

The intention of using markup forms to learn about data fields is to minimize expensive manual operations. Markup codes are typed on the areas where data are expected to be found. The field registration module locates and interprets the coded information to find field boundaries and to define data contextual constraints. The location of a markup code and its coded information are used to find the exact physical boundary of the field. The coded field type, like 'A', indicates that the field is 4-side enclosed. This helps the boundary detector to locate the boundary starting from the markup code location.

Location of a field is represented as relative distance to the form name since this distance is fixed even the absolute locations of the fields are different on different copies. This addressing scheme is less vulnerable to form displacement during image capturing. Instruction in a field is also located and "masked out" so that it will not be mistakenly considered as data in data processing stage.

6 Knowledge-based Data Extraction and Preprocessing

Extracting field data is based on the predefined field boundaries. Calculating absolute field location is based on the estimated form skew angle and the field relative distance to the form name. This adjustment scheme can avoid to de-skew the whole form prior to

data extraction. Skewed data will be extracted from the field (See Figure 5a) and then de-skewed. This approach is much more efficient since the total data areas occupy only a very portion of a form.

Extracted field data may contain noise, fragments of instruction and field separators. A data preprocessing step is required to remove non-data components from the field prior to recognition. However, occurrences of the non-data components are not necessary random. Occurrence of field separator is defined in the field markup code. Heuristic rules can be derived to clean up data of a particular field type. For example, a dotted line remover is needed for social security number fields and a slash remove is needed for date fields (See Figure 5b).

A99! 26! 9320

(a)

A99 26 9320

(b)

Figure 5: Knowledge-based data extraction and preprocessing. (a) The social security number extracted from a skewed form contains fragments of the instruction and the dotted line field separators. (b) The same field image with skewed corrected. The data can be cleaned up based on the knowledge of the instruction location and the separator structural appearance.

7 Contextual-based Data Recognition

Data interpretation is the last step and also the most important step in form processing. A knowledge-directed text understanding system is being developed which fully utilizes contextual constraints. A context-free grammar is used to encode field contextual constraints and a parser is implemented to recognize and correct recognition errors based on the syntactic definition of a field. This approach significantly simplifies the data recognition task since only a parser is needed to process different types of fields. This approach also utilizes contextual knowledge to select appropriate recognition tools for a field. For example, a home address field needs a digit recognizer to read the first word as the street number and a lexicon-based word recognizer to read street name [5]. Without applying contextual constraint prior to recognition, it is not possible to use the recognition tools and processing strategies effectively.

8 Implementation and Evaluation

The reader has been implemented in C and it runs on SUN SPARC workstations. System performance can be evaluated at two different stages. The first stage is to determine how accurate the field registration module can obtain the boundaries of the data

fields and the second stage is to evaluate text recognition capability by comparing output of the text recognizer with printed form data. In order to evaluate system performance with high confidence, it is essential to use a large variety of forms and data. The IRS form images from the National Institute of Standards and Technology [NIST] Structured Form Database provide an excellent data source for both training and testing. The database contains more than 5,500 full page form images of 12 different tax forms from the IRS 1040 Package X for the year 1988. The data on these forms are machine-printed. In addition, insurance claim forms and office correspondence forms are also used to ensure system robustness and it is not tuned to deal with tax forms only. Totally 25 different types of forms have been used in the test and 150 forms are tested.

A success rate of 95% has been achieved in field registration. 16 out of 25 forms have all the fields correctly registered. All the other forms have over 90% registration rates. A manual verification step is needed to ensure fields are corrected located or to correct system outputs.

A 100% is attained in form identification. In most cases, only topological information is needed in discrimination.

Evaluation of data extraction has also been conducted. The experiment studies the capability of data extraction using the acquired knowledge (assuming that the knowledge is complete and correct). A correct data extraction rate of 98% on non-empty data fields is attained. 123 out of 150 forms have all the data correctly extracted. Most of the errors are due to characters touching line boundaries and the system is not capable to separate the characters from the lines completely.

Finally, 94% of the field data are correctly recognized. 110 out of 150 forms have all the field data correctly recognized.

9 Summary

A robust form reader has been described. The reader is designed to handle a large variety of forms. Acquiring information about data field layout and contextual constraints of a form can be realized by presenting a sample form with markup codes to the field registration module. The field registration module detects markup codes and derives physical boundaries of the fields. Data recognition module relies on the acquired information to extract and interpret data on forms. Experiments have demonstrated that the reader can achieved high performance in field registration and data recognition. This shows that the reader is robust enough to handle a large variety of forms.

References

- [1] S.L. Taylor, R. Fritzson and J.A. Pastor, Extraction of Data from Preprinted Forms, *Machine Vision and Applications*, Vol 5, No 3, 1992, pp 211-222.
- [2] J. Yuan, L. Xu and C.Y. Suen, Form Items Extraction by Model Matching, *First International*

Conference on Document Analysis and Recognition, Saint-Malo, France, September 1991, pp 210-218.

- [3] M.M. Scholl, G.S. Himes and F.A. DeCosta, Identification and Correction of Rejection and Substitution Errors in Optical Character Recognition Systems, *IS&T/SPIE's Symposium on Electronic Imaging: Science & Technology*, San Jose, California, February 1993.
- [4] R.O. Duda and P.E. Hart, Use of the Hough Transformation to Detect Lines and Curves in Pictures, *Communications of the ACM*, Vol 15, No 1, January 1972, pp 11-15.
- [5] S.W. Lam, X. Shen and S.N. Srihari, Representing Lexicons by Modified Trie for Fast Partial String Matching, *IS&T/SPIE's Symposium on Electronic Imaging: Science & Technology*, San Jose, California, January 1993.