

# Benchmark Database and GUI Environment for Printed Arabic Text Recognition Research

AMIN G. AL-HASHIM and SABRI A. MAHMOUD

Department of Information and Computer Science

King Fahd University of Petroleum and Minerals

KFUPM 1510, Dhahran 31261

SAUDI ARABIA

e-mail: [alhasha@kfupm.edu.sa](mailto:alhasha@kfupm.edu.sa) , [smasaad@kfupm.edu.sa](mailto:smasaad@kfupm.edu.sa)

Web: <http://www.ccse.kfupm.edu.sa/~alhasha/> , <http://faculty.kfupm.edu.sa/ICS/smasaad/>

*Abstract:* - This paper presents the details of a comprehensive database of Printed Arabic text for Arabic text recognition research. It consists of scanned images of different forms of Arabic printed text (viz. book chapters, advertisements, magazines, newspapers, and reports) scanned with 200, 300, and 600 dpi resolutions. A total of 6954 pages are scanned. The database may be utilized by Arabic printed text recognition research community. It may be used as a benchmark database where researchers can evaluate their algorithms and results compared with published work of other researchers using the same database. To the best of our knowledge, there is no public comprehensive printed Arabic text database that is freely available. Hence, this database may address this deficiency in Arabic printed text recognition research. This database will be made freely available to interested researchers. In addition, this paper presents a software GUI environment to make the manipulation of the created database easier. Moreover, the software GUI provides a number of image-processing functions that can be used in the field of automatic text recognition.

*Key-Words:* - Arabic Printed Text Database, OCR Datasets, Arabic Text Recognition, OCR

## 1 Introduction

This paper presents the details of a comprehensive database of Printed Arabic text for Arabic text recognition research. It consists of scanned images of different forms of Arabic printed text (viz. book chapters, advertisements, magazines, newspapers, and reports) scanned with 200, 300, and 600 dpi resolutions. A total of 6954 pages are scanned. The database may be utilized by Arabic printed text recognition research community. It may be used as a benchmark database where researchers can evaluate their algorithms and results compared with published work of other researchers using the same database. To the best of our knowledge, there is no public comprehensive printed Arabic text database that is freely available. Hence, this database may address this deficiency in Arabic printed text recognition research. This database will be made freely available to interested researchers.

Having a public large-scale comprehensive database of Arabic text ready for researchers and developers can contribute to the development of the ATR field in different aspects. Researchers need not prepare their own data. In addition, It will enable researchers and developers to collaborate remotely and compare algorithms and results [1, 2, 3]. This paper describes a database that aims to serve Arabic

ATR researchers and developers. The database, which includes Arabic printed text, is named as PATDB (Printed Arabic Text DataBase).

Beside the PATDB, a GUI software environment is developed. This software provides two main groups of functions: database manipulation functions and image-processing functions. The first set of functions makes the manipulation of the database easier while the other set of functions helps ATR researchers by providing a number of image-processing functions that are of common use.

The rest of the paper is organized as follows: Section 2 surveys the printed Arabic text databases available in the literature. A description of the PATDB contents and attributes is given in Section 3. The process used to develop the PATDB is detailed in Section 4. Section 5 expresses the details of the PATDB and the meta information of PATDB records. The details of the software GUI environment are presented in Section 6. Finally, conclusions are presented in Section 7.

## 2 Literature Review

To the best of our knowledge, there is no public, large-scale, and comprehensive printed Arabic text database that is freely available. A small number of

bounded, commercial and/or special purpose Arabic text databases are only available. The most popular printed Arabic text databases in the literature of which the writers aware are: ERIM database and DARPA corpus.

### 2.1 ERIM Database

ERIM database was created by Environmental Research Institute of Michigan from a set of machine-printed Arabic books and magazines. It contains over 750 pages that consists of approximately 1,000,000 characters and over 200 distinct ligatures. Pages were scanned with a resolution of 300 dots per inch (dpi). The database is divided into 3 distinct sets, namely, training, statistics and testing set. It is available on a CD ROM for US\$ 40 [4].

ERIM database covers only two aspects of real life written communications, namely, books and magazines. However, many other aspects of real life written communications exist such as letters and newspapers. This coverage limitation is considered a disadvantage when developing a general-purpose ATR application. In addition, one may consider the fact that it is not freely available as another disadvantage.

### 2.2 DARPA Corpus

DARPA (Defense Advanced Research Projects Agency) Arabic corpus was created by Scientific Application International Company (SAIC) for the US Department of Defense [5]. DARPA contains 345 pages of images (around 670,000 characters) with ground-truth. Images were scanned with a resolution of 600 dpi. Images are zones of a single column of text that vary in quality. The corpus was collected from book chapters, magazine articles, newspapers and computer generated documents having only 4-fonts.

DARPA corpus has the same disadvantage of ERIM database. It does not cover all the aspect of real life written communications such as letters and advertisements. In addition, it is not currently freely available.

## 3 Database Overview

The set of document pages, that are included in the PATDB, are selected from various printing forms:

- Advertisements (2.1%),
- Book chapters (21.3%),
- Magazines (59.2%),

- Newspapers (3.0%), and
- Reports (9.9%).

Section 5 illustrates the details of each category of the selected document pages.

Based on our analysis, the needed attributes of the Arabic database are the followings:

- Simple, i.e., easy to use;
- Extensible, i.e., eligible for adding new document types without modifications;
- Standard, i.e., any addition to the database will follow a set of pre-defined rules;

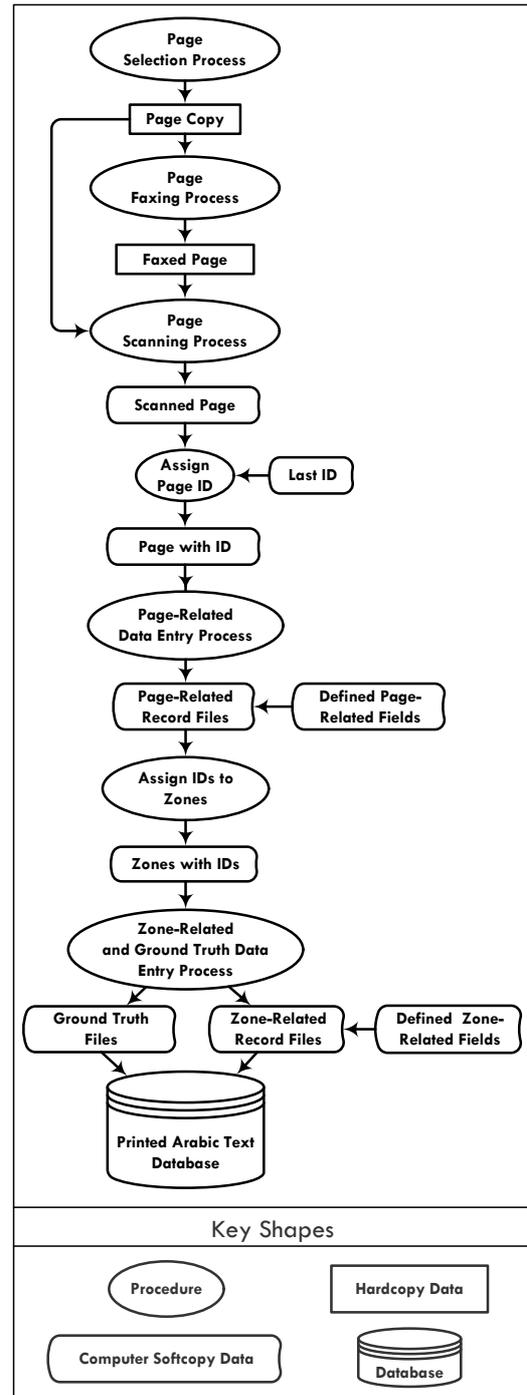


Fig.1: PATDB high-level implementation process

TABLE 1: DISTRIBUTION OF PAGE IMAGES ACROSS PATDB

Category	Color Format / Resolution									Total
	Black & White			Grayscale			Color			
	200 dpi	300 dpi	600 dpi	200 dpi	300 dpi	600 dpi	200 dpi	300 dpi	600 dpi	
ADS	22	22	14	22	22	14	22	22	14	174
BOOK	111	111	111	111	111	111	0	0	0	666
BOOK (fax)	111	187	111	187	111	111	0	0	0	818
MAG	536	536	536	536	536	536	284	284	336	4120
NEWS	35	35	35	35	35	35	0	0	0	210
REP	161	161	161	161	161	161	0	0	0	966
Total	976	1052	968	1052	976	968	306	306	350	6954

- Public;
- Comprehensive, i.e., covers different types of documents;
- Uniform, i.e., allows different sub-tasks to access different types of data in the same way;
- Reflects the physical document hierarchy as well as the logical structure;
- Stores the scanned images and their ground-truth information (corresponding text, style information, etc.) separately;
- May include a set of tools for manipulating the database and simulating the natural paper degradation models such as paper aging, multi-cropping, skewing and pepper-and-salt noise.

The construction process of the PATDB, which is addressed in Section 4, takes into account the above-mentioned attributes.

#### 4 Database Construction Process

The construction process of the PATDB follows an implementation methodology similar to the one carried out by [6]. Figure 1 outlines this implementation methodology. As shown in the figure, the process starts by selecting an appropriate page. The selection process depends on a set of pre-defined criteria. The criteria are specific to each category type. After selecting an appropriate document page, the page may be faxed before scanning to represent faxed images with lower resolution. After that, the page is scanned and given a unique number. This unique number acts as the page identifier and appears in all the files that describe this page. Then, the page is attributed and zoned. The attributes and ground-truth value of each identified zone are keyed. All the resulting files (page image, attribute files and ground-truth value files) are named according to a pre-defined naming convention, stated in Section 5.4, before uploading them to the database.

#### 5 Database Specification

This section addresses the specifications of the PATDB. The different image formats used to store the scanned page images in the PATDB are expressed in Section 5.1. Section 5.2 presents the details of the storage requirements of the PATDB. PATDB meta information is detailed in Section 5.3. Section 5.4 defines the naming conventions of the files of PATDB.

##### 5.1 Page Image Formats

The document pages, that are stored in the PATDB, are scanned and stored in three different formats:

- (1) Black & white (binary) format with color depth of 1-bit per pixel;
- (2) Grayscale format with color depth of 8-bit (1-byte) per pixel (0 to 255 gray levels); and
- (3) Color (or RGB) format with color depth of 24-bit (3-byte) per pixel.

The resolutions at which the document pages are scanned are 200, 300 and 600 dots per inch (dpi).

The uncompressed TIFF file format (file-extension .tif) is chosen to store the images of the scanned document pages in the PATDB. The TIFF format is selected because it can store complex information for the CMYK color model and can use JPEG compression techniques.

All the page images in the PATDB are scanned using HP Scanjet N8400 series scanner. The faxing of the page images is done within our campus from a low-end Panasonic fax to another Panasonic fax.

##### 5.2 Storage Requirement

The PATDB contains 6954 page images with their associated metadata files (records). These records describe the scanned pages and the zones of these pages. Three records describe the page images, namely, a page condition record, a page attribute record and a page bounding box record. Similarly, three records describe each identified zone on the

TABLE 2: PAGE CONDITION RECORD

Attribute	Possible Values
Document ID	
n-th copy	0, 1, 2, ...
n-th fax	0, 1, 2, ...
Resolution	200 (fax), 300, 600, ...
Scanning type	black & white, grayscale, color
Color depth (in bits)	1, 2, 4, 8, ...
Degradation type	original, page aging, ...
Visible salt/pepper noise	yes/no
Visible vertical streaks	yes/no
Visible horizontal streaks	yes/no
Extraneous symbols on the top	yes/no
Extraneous symbols on the bottom	yes/no
Extraneous symbols on the left	yes/no
Extraneous symbols on the right	yes/no
Page skewed on the left	yes/no
Page skewed on the right	yes/no
Page smeared on the left	yes/no
Page smeared on the right	yes/no
Visible page rotation	yes/no
Page rotation angle (in degree)	
Page rotation angle standard deviation	

page image, namely, a zone bounding box record, a zone attribute record and a zone ground-truth value record. The detailed description of these records is stated in Section 5.3. Table 1 shows the number of images available in the PATDB.

### 5.3 Meta Information

This section provides the metadata information of the page-related records (page condition record, page attribute record, and page bounding box record) and zone-related records (zone bounding box record, zone attribute record, and zone truth-value record). The metadata is adopted from [7] to reflect the PATDB.

#### 5.3.1 Page Condition Record

The purpose of the set of attributes that constitute the page condition record is to describe the various visual conditions of a given page image. These attributes and their possible values are listed in Table 2. As shown in the table, these attributes are self-explanatory. The ‘*n-th copy*’ attribute represents the number of consecutive copies applied

TABLE 3: PAGE ATTRIBUTE RECORD

Attribute	Possible Values
Document ID	
Document language	Arabic, Farsi, Dari, Azeri, Urdu, Uygur, Tajik, Pashto, Kurdish, English
Document script	Arabic, Latin
Document type	newspaper, book, report, magazine, advertisement, letter
Publication information	
Multiple pages from the same article	yes/no
Text zone present	yes/no
Displayed math zone present	yes/no
Table zone present	yes/no
Half-tone zone present	yes/no
Drawing zone present	yes/no
Page header present	yes/no
Page footer present	yes/no
Max number of text columns	
Page column layout	regular, combined-columns
Character orientation	up-left, up-right, rotated-right, rotated-left
Text reading direction	left-right, right-left, top-down, bottom-up
Dominant font type	Traditional Arabic, Arial, Simplified Arabic, Arabic Transparent, Times New Roman, Andalus, Courier New, Microsoft Sans Serif, Tahoma, ...
Dominant character spacing	proportional, fixed
Dominant font size (pts)	<< 9, 9-12, 13-18, 19-24, 25-36, >> 36
Dominant font style	plain, bold, italic, underline, other

on the page. Similarly, the ‘*n-th fax*’ attribute presents how many times the page is faxed consecutively before scanning.

#### 5.3.2 Page Attribute Record

The various properties of a given page are described through the attributes of the page condition record. The list of attributes is given in Table 3.

The value of the ‘*document language*’ attribute is ‘*Arabic*’ for PATDB. It is used for databases that may be produced in other languages later.

The ‘*publication information*’ attribute holds the information about the source from which the page is taken.

The ‘*multiple pages from the same article*’ attribute indicates whether an article span more than one page. This attribute can be used with the ‘*publication information*’ attribute to retrieve a full article that exists on more than one page.

The ‘*max number of text columns*’ represents the number of text columns in the body area of the given page.

The ‘*character orientation*’ attribute provide the orientation of the characters within the text line when the page oriented to up-right position. Similarly, the ‘*text reading direction*’ attributes presents the text reading direction within a text line when the page is oriented to up-right position.

### 5.3.3 Page Bounding Box Record

The boundary of the header area, body area and footer area of a given document page is defined in the page bounding box record. For each area, the coordinates of the upper-left and lower-right corners are included in the record. Table 4 shows the six attributes of these three areas.

TABLE 4: PAGE BOUNDING BOX RECORD

Attribute	Possible Values
Document ID	
Header upper-left corner coordinates	(X,Y)
Header lower-right corner coordinates	(X,Y)
Live matter area upper-left corner coordinates	(X,Y)
Live matter area lower-right corner coordinates	(X,Y)
Footer upper-left corner coordinates	(X,Y)
Footer lower-right corner coordinates	(X,Y)

### 5.3.4 Zone Bounding Box Record

The zone bounding box record of a given zone in a given document page contains its identification number and its upper-left and bottom-right bounding box (x, y) coordinates within the document page.

### 5.3.5 Zone Attribute Record

The zone attribute record describes the common characteristics of an identified zone in a given document page. Table 5 presents the set of attributes that constitute the zone attribute record.

The ‘*zone’s column number*’ attribute describes the zone’s column location. A zone may be in the header area, footer area and column 1 of 1, 1 of 2, etc.

TABLE 5: ZONE ATTRIBUTE RECORD

Attribute	Possible Values
Document ID	
Zone ID	
Zone content	text, text with special symbols, displayed math, table, half-tone, drawing, form, ruling, bounding box, logo, map, advertisement, announcement, handwriting, others
Text zone label	text body, list item, drop cap, caption, abstract body, abstract heading, section heading, synopsis, highlight, pseudo-codes, reference heading, reference list item, footnote, author biography, page header, page footer, page number, article title, author, affiliation, diploma information, society membership information, article submission information, abstract heading, abstract body, footnote heading, keyword heading, keyword body, other
Text alignment within the zone	left aligned, center aligned, right aligned, justified, justified hanging, left hanging
Dominant font type	Traditional Arabic, Arial, Simplified Arabic, Arabic Transparent, Times New Roman, Andalus, Courier New, Microsoft Sans Serif, Tahoma
Dominant character spacing	proportional, fixed
Dominant font size (pts)	<< 9, 9-12, 13-18, 19-24, 25-36, >> 36
Dominant font style	plain, bold, italic, underline, other
Character orientation	up-left, up-right, rotated-right, rotated-left
Text reading direction	left-right, right-left, top-down, bottom-up
Zone’s column number	
Next zone ID within the same thread	

The zones of each document page can be grouped into several logical units. Within each logical unit, the reading order is sequential. This logical unit is called a semantic thread. So, the ‘*next zone ID within the same thread*’ attribute is used to indicate the reading order among the zones that

constitute a semantic thread. 'nil' is used to indicate the end of the semantic thread.

### 5.3.6 Zone Truth-Value Record

The zone truth-value record is available in the case of text zones only. It contains the ground-truth value of the given text zone.

All the page-related and zone-related records, except the zone truth-value record, are stored as ASCII text files (extension .txt) with a simple representation using JSON (JavaScript Object Notation). JSON bases on name/value pairs idea [8]. In PATDB, the name represents the attribute name while the value is one of the possible values of this attribute.

The zone truth-value records are UNICODE text files (extension .txt) that holds only the ground truth-value of a given text zone.

A part of a text file represented using JSON for a page attribute record of a given book page image is given as an example in Figure 2.

```
{
  "document id" : "BOK0001P001_C00F00_R0200B",
  "document language" : "Arabic",
  "document script" : "Arabic",
  "document type" : "book",
  "publication information" : "Alaalami Library,Beirut",
  "multiple pages from the same article" : true,
  "text zone present" : true,
  "special symbol present in text zone" : false
}
```

Fig.2: A Sample of a Page Attribute Record File

### 5.4 File Naming and Database Hierarchy

Table 6 and Table 7 show the naming conventions of the different categories of the scanned document pages and the scanning types, respectively.

TABLE 6: CATEGORIES ABBREVIATIONS

Abbreviated Name	Full Name
ADS	Advertisements
BOOK	Book Chapters
MAG	Magazines
NEWS	Newspapers
REP	Reports

TABLE 7: SCAN-TYPE ABBREVIATIONS

Abbr.	Meaning	Color Depth (bit/pixel)	Number of Colors
BW	Black & White	1	2
GS	Grayscale	8	256 gray shades
CL	Color	24	Millions of colors

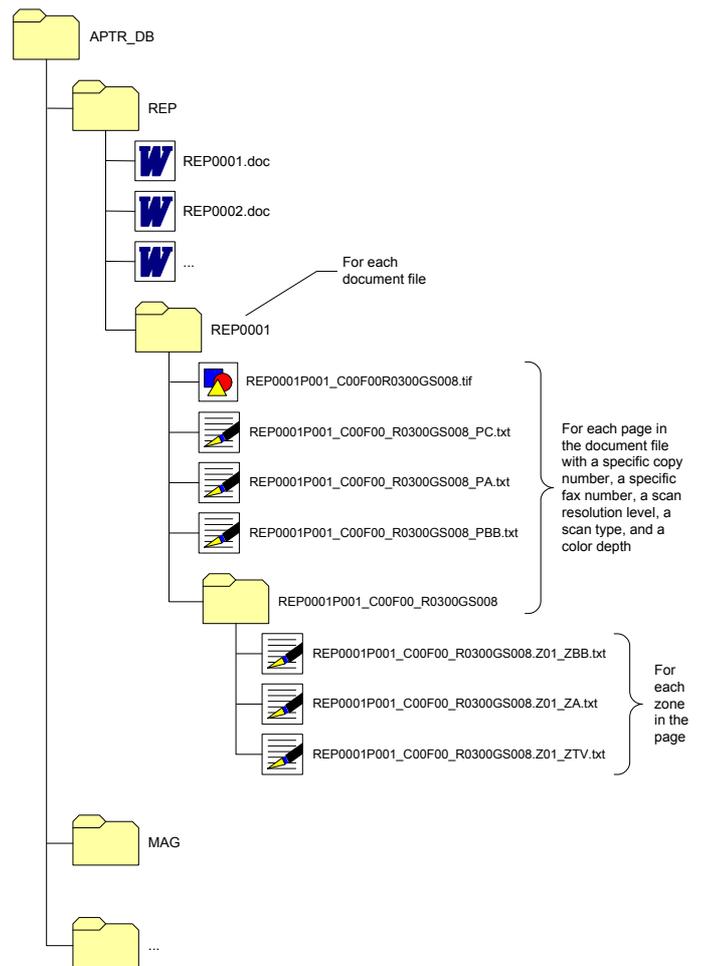


Fig.3: PATDB Hierarchy

All the items (page images and records, etc.) that make up the PATDB are placed and named according to the following criteria:

- The first level of the database contains a list of directories. Each one of these directories represents a category of document files, e.g., magazines. These directories are named according to the 'Categories Abbreviations' table. Therefore, any newly acquired document file should be added to one of these categories. If the document file cannot fit to any of these directories, then a new directory is created after recording its name and abbreviation in the 'Categories Abbreviations' table.
- The name of the document file within its corresponding directory is formatted as follows: *CategoryNamennn.FileExtension*, where *CategoryName* is the abbreviated category name according to the 'Categories Abbreviations' table; *nnn* is a 4-digit unique number; and *FileExtension* is the extension of the file being stored. This document file acts as

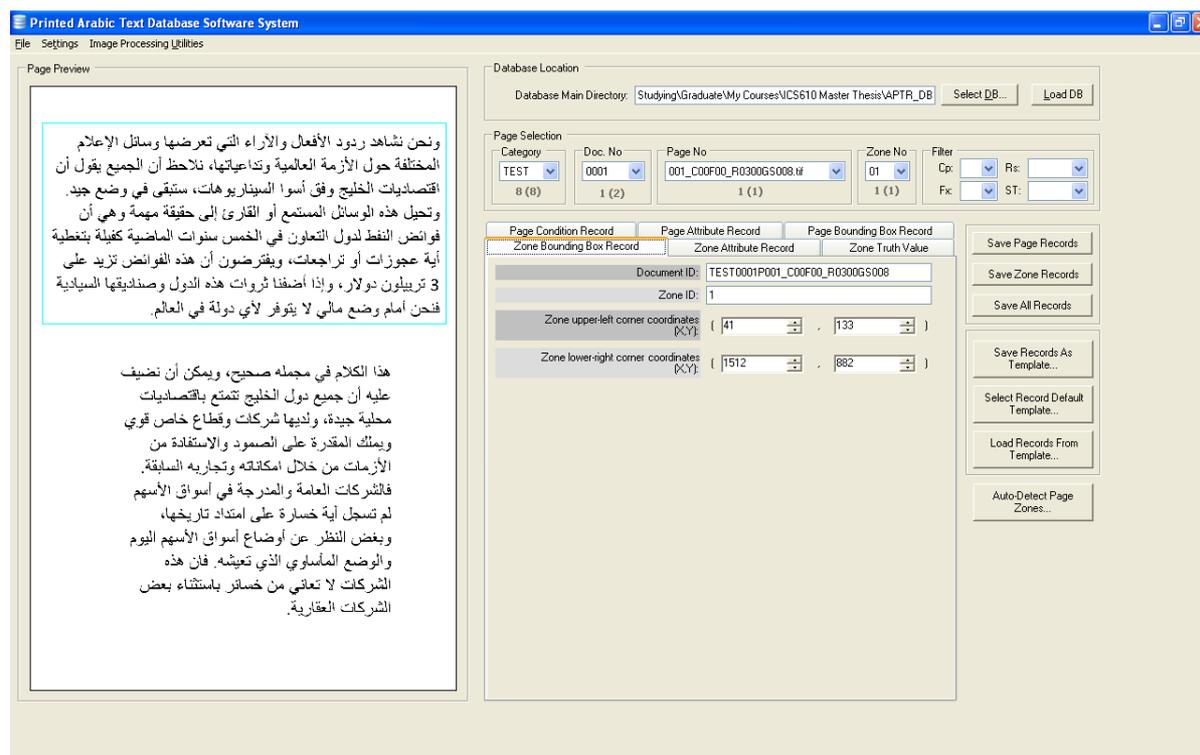


Fig.4: A snapshot of the main form of the software system

the source from which the ground-truth value of the images is taken.

- For each added document file, a directory is named after the document file name excluding the file extension. This directory is created in the corresponding category directory. It includes the following items for each page of the document file:

1. A scanned image is named as follows:  
*DocNamePppp\_CccFff\_RrrrrScanTypeIll.tif*,  
where DocName is the document name without file extension; ppp is a 3-digit page number; cc is a 2-digit copy number; ff is a 2-digit fax number; rrrr is a 4-digit resolution level in dot per inch (dpi); ScanType is an abbreviated scan type according to the 'scan types abbreviations' table; and Ill is a 3-digit color depth in bits.  
The copy number, referred to by cc, represents the number of consecutive copies that are done on the page before scanning it. Similarly, the fax number, referred to by ff, represents the number of consecutive faxes that are done on the page before the scanning process.
2. Three text files representing the page condition record, page attribute record, and page bounding box record are named as follows, respectively:  
*DocNamePppp\_CccFff\_RrrrrScanTypeIll\_PC.txt*,  
*DocNamePppp\_CccFff\_RrrrrScanTypeIll\_PA.txt*, and  
*DocNamePppp\_CccFff\_RrrrrScanTypeIll\_PBB.txt*,

where DocName, ppp, cc, ff, rrrr, ScanType and Ill are as defined earlier.

3. A directory for the page zones is named as follows:

*DocNamePppp\_CccFff\_RrrrrScanTypeIll*,  
where DocName, ppp, cc, ff, rrrr, ScanType and Ill are as defined earlier.

This directory includes three text files for each identified zone in the page. These text files represent the zone bounding box record, zone attributes record and zone ground-truth value record and they are named as follows, respectively:

*DocNamePppp\_CccFff\_RrrrrScanTypeIll.Zzz\_ZBB.txt*,  
*DocNamePppp\_CccFff\_RrrrrScanTypeIll.Zzz\_ZA.txt*, and  
*DocNamePppp\_CccFff\_RrrrrScanTypeIll.Zzz\_ZTV.txt*,  
where DocName, ppp, cc, ff, rrrr, ScanType and Ill are as defined earlier and zz is a 2-digit zone number.

Figure 3 presents the above representation graphically. The figure shows in details a sample case from the report category (REP). The other types are similarly constructed.

## 6 GUI Software Environment

Along with the PATDB, a software system is developed (Figure 4). This software system works

as a management system for the PATDB. The sole purpose of this software system is to make the job of those who would like to upload more documents to the database easier. It enables the users to upload more documents and manipulate the PATDB through an easy-to-use graphical user interface (GUI) forms. In addition, the software system provides the ATR researchers and developers with a number of basic image processing utilities that may be used in their field.

The software system provides two broad sets of functions: database-related functions and ATR-related functions. The ATR-related functions can be reached mainly through the 'Image Processing Utilities' menu while the rest of menus server the database-related functions. Each set of these functions is discussed in the following sub-sections.

## 6.1 Database-Related Functions

The software system offers a set of functions through which the user can browse and manipulate any database following the file naming conventions of PATDB specification shown in Section 5. This set of functions enables the user to:

- Traverse through the available document pages one by one with their associated description files (page condition record, page attribute record and page bounding box record) shown aside separately.
- Traverse through the available zones of any selected document page with the zones associated description files (zone bounding box record, zone attribute record and zone ground-truth value record) presented aside separately.
- Modify any description file associated with any document page or zone through easy-to-use forms.
- Save the description files associated with any document page or zone as template and then load these saved templates to other document pages and zones.

## 6.2 ATR-Related Functions

The software system can assist the ATR systems by providing a set of ATR-related functions. These functions can be used in the early stages of the processing process of an ATR system in order to enhance the accuracy level. Following sub-sections explain each one of these implemented functions.

### 6.2.1 Bounding Box Detector Function

The bounding box of a zone (image, figure, text, etc.) on a page is the rectangle that surrounds it ignoring the surrounding space. The bounding box is described by four numbers: the x-y coordinates of the upper-left corner and the x-y coordinates of the lower-right corner of the zone. The coordinates are measured, in pixels, from the top-left corner of the page. The user can modify the detected bounding box according to his preference.

The bounding box detector function provided by the software system is ideal in detecting the bounding boxes of zones on pages with single-column layout. Figure 5 shows the bounding box rectangles detected by the function of the two available text paragraphs on a sample single-column page. In addition, it shows the coordinates of these two bounding boxes in the 'Detected Zones' area.

#### Zone bounding box identification Algorithm

Detecting the bounding box rectangle of a zone (a text paragraph or a text line) on a single-column page involves the following steps (adopted from [9]):

##### 1. Horizontal Projection

Project the page horizontally pixel by pixel. Then, store the result into an array,  $hArray$ , of size equals to the page height in pixel.

##### 2. Bounding Box Y-Coordinates

Find the bounding box y-coordinates of the zone by conducting the following two steps:

###### • Step 1 (Top Y-Coordinate)

Traverse forward through  $hArray$  to find the y-coordinate of the upper-left corner of the bounding box,  $y^{top}$ . If the zone is the top most one in the page, then  $y^{top}$  is the zero-value index directly before the first non-zero-value index of  $hArray$ . The other zones  $y^{top}$  is calculated in the same manner except that the first non-zero-value index must be preceded by at least  $n$  consecutive zero-value indices. These  $n$  zero-value indices, which are determined by the user, represent the minimum vertical spaces (in pixel) that must separate any two consecutive zones. By having a small value for  $n$ , the algorithm detects the non-intersecting text lines of the pages. Similarly, the bounding boxes for each paragraph in the page can be detected if we increase the value of  $n$ . Likewise, the bounding box of the whole page content can be detected by setting  $n$  to a large number or the page height.

###### • Step 2 (Bottom Y-Coordinate)

Traverse through the  $hArray$  again starting from  $y^{top}$  until a zero-value index is reached. This zero-value index represents the y-coordinate of the lower-right corner of the zone bounding box rectangle,  $y^{bottom}$ .

##### 3. Bounding Box X-Coordinates

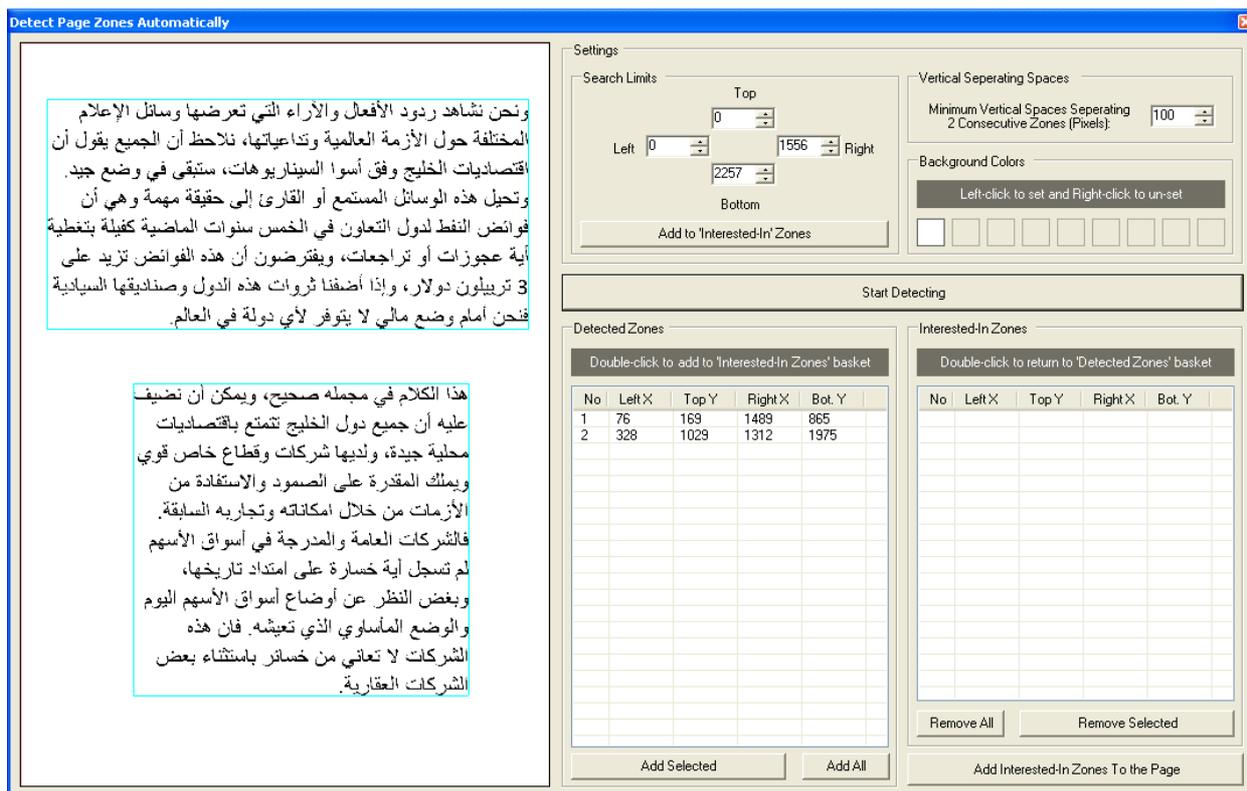


Fig.5: A snapshot of the bounding box detector function carried out on a sample document image

Find the bounding box x-coordinates of the zone by carrying out the following three steps:

- **Step 1 (Vertical Projection)**

Project the page vertically pixel by pixel starting from  $y^{top}$  to  $y^{bottom}$  and store the result into an array,  $vArray$ .

- **Step 2 (Left X-Coordinate)**

Traverse forward through  $vArray$  until a non-zero-value index is reached. The zero-value index directly preceding this non-zero-value index represents the x-coordinate of the upper-left corner of the bounding box rectangle,  $x^{left}$ .

- **Step 3 (Right X-Coordinate)**

Traverse backward through  $vArray$  until a non-zero-value index is reached. The zero-value index directly preceding this non-zero-value index is the x-coordinate of the lower-right corner of the bounding box rectangle,  $x^{right}$ .

4. Repeat steps 2 and 3 until the end of the  $hArray$  is reached.

### 6.2.2 De-Noising Functions

Any unwanted information that comes out into images from any source (e.g., data acquisition process, naturally occurring phenomena, etc.) is considered noise [10, 11]. For example, when a document page is scanned and then digitized, a certain amount of noise may appear. This noise

incorporates difficulties to the recognition process of the ATR systems [3, 12]. It may lower the recognition accuracy level of the system and even cause failure if the system was not developed to address this kind of noise [13]. Different approaches can help in this regards. They can help in removing the incorporated noise from the scanned image before processing it by the ATR system. Following sections presents two algorithms for de-noising implemented by the software system.

#### 6.2.2.1 Statistical Based Smoothing Function

The statistical based smoothing function tackles mainly the noise pixels that add irregularities to the outer boundary of the characters. This function reduces the incorporated noise on the binary images by getting rid of small areas and filling little holes that make the character contour regular [14]. Filling and deletion depending on the pixel's initial value and its neighbors' initial values. The function can handle only black-&-while (binary) images.

#### Algorithm

The algorithm of the statistical based smoothing function bases on a statistical decision criterion. Given a binary image, the function modifies (fills or eliminates) each pixel depending on the pixel's initial value and its neighbors' initial values. The rules, taken from [8], are stated as follows:

If  $P_0 = 0$  then

$$P'_0 = \begin{cases} 0, & \text{if } \sum_{i=1}^8 P_i < T \\ 1, & \text{otherwise} \end{cases}$$

else

$$P'_0 = \begin{cases} 1, & \text{if } P_i + P_{i+1} = 2 \text{ for at least one } i = 1, \dots, 8 \\ 0, & \text{otherwise} \end{cases}$$

where  $P_0$  is the current pixel value,  $P'_0$  is the new pixel value and  $T$  is the threshold. The zero '0' in the above rules means white pixel while the one '1' means black pixel. The labeling scheme of these pixels is shown in Figure 6.

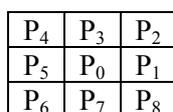


Fig.6: The current pixel  $P_0$  and its neighbors

### 6.2.2.2 Average Smoothing Function

The average smoothing function attempts to smooth the image edges and corners by filling small holes or deleting small fills. The filling and deletion is determined by a  $3 \times 3$  weighted matrix elements. The system gives defaults parameters. However, the user can modify the parameters. The function can handle both black-&-white (binary) and grayscale images.

#### Smoothing Algorithm

Given a  $3 \times 3$  weighted matrix elements of a total sum equals to one (see Figure 7-C), the steps of the average smoothing function (adopted from [12]) are as follows:

1. Scan the image pixel by pixel with a window size of  $3 \times 3$  as shown in Figure 7-B.
2. Multiply each pixel in the  $3 \times 3$  image window by its corresponding element in the  $3 \times 3$  weighted matrix elements. Then, sum all of the multiplications into variables called *total*. Figure 7-D shows the result of multiplying a sample  $3 \times 3$  image window by the weighted matrix elements.
3. Round the total variable into its nearest integer value as shown in Figure 7-E.
4. Change the value of the pixel at the center of the  $3 \times 3$  image window to the new pixel that has its value equals to the value of the *total* variable.

### 6.2.3 Page De-Skew Function

When a document page is copied or scanned, a slight degree of skew may be introduced. In order to achieve good recognition results, this skew need to

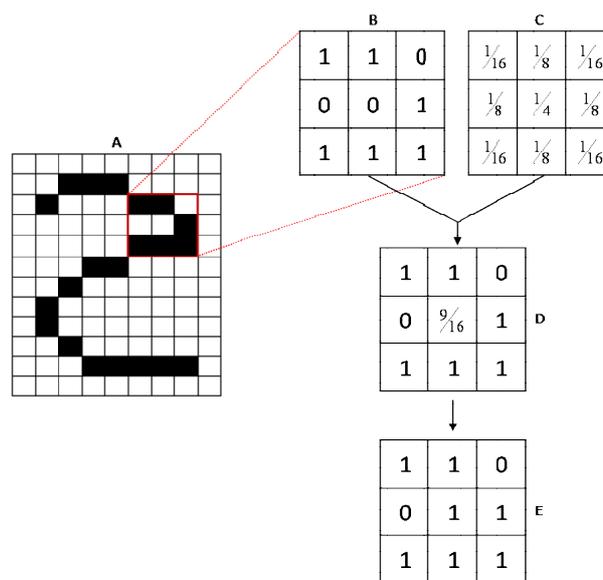


Fig.7: A sample iteration of the average smoothing function on a sample black-&-white page image (A). (B) is the initial  $3 \times 3$  image window. (C) is the  $3 \times 3$  weighted matrix elements. (D) is (B) after applying Step 2. (E) is (D) after applying Step 3 then Step 4.

be corrected before passing the page image into the ATR system. The aim of this function is to enable the user to correct such skew through few mouse-clicks.

#### De-skew Algorithm

Given a page image that contained skewed text, the function de-skews the text by carrying out the following steps:

1. The user is asked to draw a line that simulates the base line of the dominant text on the given page.
2. The function calculates the slope absolute value of the drawn line and finds the slope angle, *angle*, in degree.
3. The function rotates the page *angle* degree.

## 7 Conclusion

The details of a comprehensive database of Arabic printed text are presented. It consists of scanned images of book chapters, advertisements, magazines, newspapers, and reports that are scanned at 200, 300, and 600 dpi resolutions. A total of 6954 pages are scanned. The database may be utilized by Arabic printed text recognition research community. Thus providing researchers and developers working on the field of automatic Arabic text recognition with benchmark database of printed Arabic text which enables them to collaborate remotely and compare algorithms and results. The database will soon be made freely available for researchers and

developers. The database covers most of the written communications faced in real life. This database addresses one of the obstacles of conducting automatic Arabic printed text recognition research as to authors knowledge no such database is freely available to researchers. In addition, a GUI software environment is developed to make the manipulation job of the database easier. The software system also implements a number of image-processing functions that can help the automatic text recognition researchers in their work.

We would like to note that this database will be made freely available to interested researchers.

*References:*

- [1] J. J. Hull, *A Database for Handwritten Text Recognition Research*, IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol.16, No.5, 1994, pp. 550-554.
- [2] V. Märgner and M. Pechwitz, *Synthetic Data for Arabic OCR System Development*, ICDAR, Proceedings of the Sixth International Conference on Document Analysis and Recognition, 2001, pp. 1159-1163.
- [3] N. B. Amor and N. E. B. Amara, *A hybrid approach for Multifont Arabic Characters Recognition*, Proceedings of the 5th WSEAS Int. Conf. on Artificial Intelligence, Knowledge Engineering and Data Bases, 2006, pp. 194-198.
- [4] [www.irim.org](http://www.irim.org) (22-Jan-2008)
- [5] R. Davidson and R. Hopely, *Arabic and Persian OCR Training and Test Data Sets*, Proceedings of Symposium. on Document Image Understanding Technology, 1997.
- [6] I. T. Phillips et al., *The Implementation Methodology for CD-ROM English Document Database*, Proceedings of the Second International Conference on Document Analysis and Recognition, 1993, pp. 484 - 487.
- [7] I. T. Phillips, S. Chen and R. M. Haralick, *CD-ROM Document Database Standard*, Proceedings of the Second International Conference on Document Analysis and Recognition, 1993, pp. 478-483.
- [8] Introducing JSON. <http://www.json.org/> (02-Feb-2008)
- [9] S. Brook and Z. Al Aghbari, *Classification of Personal Arabic Handwritten Documents*, WSEAS Transactions on Information Science & Applications, Vol.5, No.6, 2008, pp. 1021-1030.
- [10] I. Rosca, L. State, and C. L. Cocianu, *Learning Schemes in Using PCA Neural Networks for Image Restoration Purposes*, WSEAS Transactions on Information Science and Applications, Vol.5, No.7, 2008, pp. 1149-1159.
- [11] V. Niola and G. Quaremba, *Image denoising and Fuzziness Measures*, Proceedings of the 7th WSEAS International Conference on Neural Networks, 2006, pp. 1-6.
- [12] C. Boiangiu and A. Dvornic, *Methods of Bitonal Image Conversion for Modern and Classic Documents*, WSEAS Transactions on Computers, Vol.7, No.7, 2008, pp. 1081-1090.
- [13] C. Boiangiu and A. Dvornic, *Bitonal Image Creation for Automatic Content Conversion*, 9th WSEAS International Conference on Automation and Information, 2008, pp. 454-459.
- [14] S. A. Mahmoud, *Arabic Character Recognition Using Fourier Descriptors and Character Contour Encoding*, Pattern Recognition, Vol.27, No.6, 1994, pp. 815-824.