

INTEGRATION OF OPTICAL CHARACTER RECOGNITION (OCR) TECHNOLOGY WITH A DATABASE APPLICATION



By:

**Annalie Faith P. Abella
Geraldine Gail A. Callos
Chryshyll T. Ti**

**SCHOOL OF ARTS AND SCIENCES
ATENEO DE DAVAO UNIVERSITY**

MARCH 2006

**INTEGRATION OF OPTICAL CHARACTER RECOGNITION (OCR)
TECHNOLOGY WITH A DATABASE APPLICATION**

An Independent Research

Presented To

The Faculty of the Computer Studies Division

Ateneo de Davao University

In Partial Fulfillment

of the Requirements for the Degree

Bachelor of Science Major in Information Technology

by

Abella, Annalie Faith P.

Callos, Geraldine Gail A.

Ti, Chryshyll T.

SCHOOL OF ARTS AND SCIENCES

ATENELO DE DAVAO UNIVERSITY

MARCH 2006

TABLE OF CONTENTS

RECOMMENDATION FOR ORAL DEFENSE.....	iii
RECOMMENDATION FOR ACCEPTANCE.....	iv
ACKNOWLEDGEMENTS.....	v
TABLE OF CONTENTS.....	vii
LIST OF FIGURES.....	xi
LIST OF TABLES.....	xii
ABSTRACT.....	1
INTRODUCTION.....	2
1.1 Background of the Study.....	2
1.2 Application Context.....	3
1.3 Objectives of the Study.....	4
1.4 Scope and Limitation of the Study.....	4
1.4.1 Scope.....	4
1.4.2 Limitations.....	5
1.5 Significance of the Study.....	5
1.6 Glossary of Terms.....	6
REVIEW OF RELATED WORKS.....	8
2.1 Software Development Toolkits.....	8
2.1.1 <i>LeadTools OCR Programming Tools</i>	9
2.1.2 <i>Asprise OCR Software Development Kit</i>	9
2.1.3 <i>SmartScan Xpress Software Development Kit</i>	9
2.1.4 <i>XP IDEA OCR Software Development Kit</i>	10
2.2 Review of Software Development Kits.....	10
2.2.1 <i>Program/Application Deficiency</i>	10
2.2.2 <i>Recognition Technologies Offered</i>	11
2.2.3 <i>Software Development Kits Information</i>	12
2.2.4 <i>Criticism of Software Development Kits</i>	12
2.3 Xceed Components for .NET.....	16
2.4 Recognizing Technologies.....	16
2.4.1 <i>Optical Mark Recognition (OMR) Technology</i>	17
2.4.2 <i>Intelligent Character Recognition (ICR) Technology</i>	17
2.4.2.1 <i>ICR Recognition Engines</i>	18
2.4.3 <i>Magnetic Ink Character Recognition (MICR) Technology</i>	19
METHODOLOGY.....	20
3.1 Survey of OCR Software Development Kit.....	20

3.2 Deciding on which OCR Software Development Kit to be used.....	20
3.3 Designing a Prototype.....	20
3.4 Development of the Application.....	21
3.5 Reviewing and Checking of the Application.....	21
3.6 Consultation with the Advisers.....	22
3.7 Finalizing the Design and Implementation of the Application.....	22
3.8 User Acceptance Test.....	22
TECHNOLOGICAL BACKGROUND.....	23
4.1 Definition of Optical Character Recognition.....	23
4.2 Optical Character Recognition Limitations.....	24
4.3 Optical Character Recognition (OCR) Engines.....	26
4.3.1 MOR OCR Engine.....	26
4.3.2 MTX (MText) OCR Engine.....	27
4.3.3 FireworX OCR Engine.....	28
4.4 Definition of Database.....	28
4.5 Relation Database Model.....	29
4.6 Components.....	30
4.7 Operational Framework.....	31
RESULTS AND DISCUSSIONS.....	33
5.1 Integrated Optical Character Recognition (OCR) Database Application....	33
5.2 Integration of Optical Character Recognition Technology with a Database Application.....	33
5.3 LeadTools Programming Tool Functions.....	35
5.3.1 Scanning.....	35
5.3.2 Zoning.....	35
5.3.3 Saving and Opening Image Files.....	36
5.3.4 Viewing.....	37
5.3.5 Paging.....	37
5.3.6 Recognizing Document Pages.....	38
5.4 Conducted a Qualitative Survey.....	39
5.5 Project Resources.....	43
5.6 Program Modules.....	44
5.6.1 Recognize Modules.....	44
5.6.2 Storing to Database Module.....	44
5.6.3 Zoning Modules.....	45
5.6.3.1 AddZone Module.....	46
5.6.3.2 DeleteZone Module.....	46
5.6.3.3 UpdateZone Module.....	47
5.6.4 Templates Module.....	47

5.6.4.1 <i>SaveTemplates Module</i>	47
5.6.4.2 <i>EditTemplates Module</i>	47
5.6.4.3 <i>DeleteTemplates Module</i>	47
5.6.4.4 <i>OpenTemplates Module</i>	48
5.6.4.5 <i>ViewTemplates Moudle</i>	48
5.6.4.6 <i>CloseTemplates Module</i>	48
5.6.5 <i>Page Module</i>	50
5.6.5.1 <i>InsertPage Module</i>	50
5.6.5.2 <i>RemovePage Module</i>	50
5.6.6 <i>Acquire Image Module</i>	50
5.6.6.1 <i>Open Image Module</i>	50
5.6.6.2 <i>Scan Image Module</i>	51
5.7 <i>Program Limitations</i>	52
CONCLUSION AND RECOMMENDATIONS	53
6.1 <i>Conclusion</i>	53
6.2 <i>Recommendations</i>	53
BIBLIOGRAPHY	55
APPENDICES	57
APPENDIX A User Interface	58
A.1 <i>Login Form</i>	58
A.2 <i>Main Form</i>	59
A.3 <i>Zone Form</i>	60
A.4 <i>Insert Page Form</i>	61
A.5 <i>Open Template Form</i>	61
A.6 <i>Select Source Form</i>	62
A.7 <i>Scan Image Form</i>	62
A.8 <i>Save Image Form</i>	63
A.9 <i>Check Database Fields Form</i>	63
A.10 <i>Compare Current Database Fields with Template Fields Form</i>	64
A.11 <i>Store to Database Form</i>	65
A.12 <i>View Data Form</i>	66
APPENDIX B. Source Codes	67
B.1 <i>Image Methods</i>	67
B.1.1 <i>Select Source</i>	67
B.1.2 <i>Acquire Image</i>	67
B.1.3 <i>Start OCR Engine</i>	68
B.1.4 <i>Shut down OCR Engine</i>	69
B.1.5 <i>Save Image</i>	69
B.2 <i>Zoning Methods</i>	70

B.2.1 Drawing/mapping zones.....	70
B.2.2 Add Zones.....	77
B.2.3 Delete Zones.....	83
B.2.4 Update Zones.....	83
B.2.5 Look-up Zone Names.....	84
B.2.6 Assign Zone Section Names.....	85
B.3 Page Processing Methods.....	86
B.3.1 Flip Image.....	86
B.3.2 Rotate Image.....	87
B.3.3 Skew Image.....	87
B.3.4 Auto-Orient Image.....	88
B.3.5 Zoom In Image.....	89
B.3.6 Zoom Out Image.....	89
B.3.7 Insert Page.....	90
B.3.8 Remove Page.....	92
B.3.9 Show Page.....	94
B.3.10 Update View Page.....	94
B.4 Processing Recognized Data.....	95
B.4.1 Recognizing Scanned Data.....	95
B.4.2 Assigning to Database Fields.....	97
B.4.3 Storing to Database.....	99
B.5 Template Methods.....	102
B.5.1 Save Template.....	102
B.5.2 EditTemplate.....	106
B.5.3 Delete Template.....	108
B.5.4 Load Template Details.....	109
B.6 Database-Related Methods.....	112
B.6.1 Checking of Database Fields Method.....	112
B.6.2 Save Database Fields Method.....	113
APPENDIX C. Database Stored Procedures.....	115
C.1 Retrieve Information Queries.....	115
C.2 Adding, Deleting and Updating Queries.....	120
APPENDIX D. Sample Questionnaire.....	124
APPENDIX E. Sample Forms.....	126
APPENDIX F. Installation Instructions.....	133
CURRICULUM VITAE.....	135
Curriculum Vitae A. Annalie Faith P. Abella.....	136
Curriculum Vitae B. Geraldine Gail A. Callos.....	137
Curriculum Vitae B. Chryshyll T. Ti.....	138

LIST OF FIGURES

Figure 4.7	Operational Framework.....	22
Figure 5.6.2	Store To Database.....	45
Figure 5.6.3	Zone Properties.....	46
Figure 5.6.4.2	Edit Templates.....	48
Figure 5.6.4.4	Open Templates.....	49
Figure 5.6.4.5	View Templates.....	49
Figure 5.6.5.1	Insert Page.....	50
Figure 5.6.6.1	Open Image.....	51
Figure 5.6.5.1	Scan Image.....	52

ABSTRACT

Nowadays, everything must be finished before or within the deadlines. People are doing so many things in so little time. Developers see this trend and decided to focus on building applications that will make the people's works easier. They began introducing and continuously developing new gadgets that can help the society finish their tasks in a quicker manner. In these past years, employees are spending most of their time manually inputting important information from different forms to the database. There are new technologies that are available and useful to these employees but they are not fully utilized by them. These technologies, like the Optical Character Recognition, can help enhance the efficiency on working environments. This document discusses the proponents' desire to expand the usage of Optical Character Recognition so that there will be no need to manually enter all information from a document individually. It discusses the integration of an optical character recognition to an application database which will lead to the efficiency of storing data electronically.

Keywords:

Optical character recognition, Application database

CHAPTER I

INTRODUCTION

1.1 Background of the Study

Over the years, database applications have been extensively used by companies in storing and acquiring accurate data in an efficient and effective manner. It helps companies organize their data to reduce inconsistency and redundancy. In addition, it also provides security restrictions causing information to be inaccessible to unauthorized people. At present, business establishments are relying more on technology in handling their business transactions. Database applications are used to reduce physical storage of information by allowing users to enter data into the database. However, database management system users still perceive a need to continuously improve the role of database system applications in a company's working environment.

Besides database applications, there are other technologies that can improve data entry and maintenance such as optical character recognition. The use of Optical Character Recognition is one of the thousands of software developed to help people in this time of rapid execution. Optical character recognition (OCR) is the translation of optically scanned bitmaps of printed or written text characters into character codes that can later be used by the computer. This is an efficient way to turn hard-copy materials into data files that can be edited and otherwise manipulated on a computer. Suppose you want to digitize a certain book overnight. You could stay up all night typing and still not

finish or you could use a high-end scanner and in minutes scan the book into a computer using optical character recognition (OCR) technology. This is the technology long used by libraries and government agencies to make lengthy documents quickly available electronically. Nowadays, OCR technologies have proven their usefulness in efficient data storage.

1.2 Application Context

Even though Optical Character Recognition and Database Applications prove to be extremely useful to its users, they are still not fully utilized by users and developers. Since developers and users want efficient data storage and retrieval, the proponents came up with the idea of integrating the two technologies to achieve this goal. This leads the proponents' study to focus on building an application that integrates Optical Character Recognition into a database application. The study's focal point will be storing the required data or information details into the database from a source document or an image file using the coordinates specified by the user. The user scans the document and then the application will then recognize the data and store the specific document details into the database.

1.3 Objective of the Study

The study aims to attain the following project objectives:

- Expand or develop the usage of a database applications by integrating it with an Optical Character Recognition technology that can make it more functional to its users
 - The technology-based application aims to integrate the Optical Character Recognition technology into a database application
 - The technology-based application aims to store accurate recognized characters into the database
- The proponents aim to familiarize and understand strategies and methods used in the OCR technology's text recognition processes that are utilized in the Optical Character Recognition Software Development Toolkits.

1.4 Scope and Limitation of the Study

1.4.1 Scope

The study will primarily focus on integrating an optical character recognition (OCR) into a specific database application by storing specific data taken from the recognized image.

1.4.2 Limitations

The study and application will not focus on resolving the existing database and optical character recognition technology's limitations. The study will simply focus on utilizing their capabilities by developing an application that would be able to recognize and store specific texts from an image file to a database.

1.5 Significance of the Study

This study is significant because it expands the usage and capabilities of the database application. It is significant because using the OCR technology in entering data into the database would hasten the encoding process. Through integrating optical character recognition into a database application, the user will be able to directly store specific data to a database either from a saved image file or an image file created after scanning an original physical document. Through this application, the process of storing data to the database will be faster because the user will not need to encode all of the information needed.

This application is useful for offices that handle various forms that need to be stored electronically. An example of these forms and offices are registration forms in the Registrar's office.

Lastly, this study can serve as a stepping stone in applying an Optical Character Recognition technology with a database application. It will not only emphasize a new way of storing data but it can open new doors for other

developers in enhancing the application's capabilities and resolving some limitations of the project.

1.6 Glossary of Terms

- 1. Optical Scanner** – is a device that can read text or illustrations printed on paper and translate the information into a form the computer can use.
- 2. ICR (Intelligent Character Recognition)** - is the computer translation of hand printed and written characters. Data is entered from hand-printed forms through a scanner, and the image of the captured data is then analyzed and translated by sophisticated ICR software.
- 3. OMR (Optical Mark Recognition)** – is the technology of electronically extracting intended data from marked fields, such as checkboxes and fill-infields, on printed forms
- 4. MICR (Magnetic Ink Character Recognition)** – is a character recognition system that uses special ink and characters. When a document that contains this ink needs to be read, it passes through a machine, which magnetizes the ink and then translates the magnetic information into characters.
- 5. ASCII Code (American Standard Code for Information Interchange)** – Assigns a number to each character that can be typed on a computer keyboard. The code aids programming and communication between computers.

6. **TWAIN** - defines a standard software protocol and application programming interface (API) for communication between software applications and image acquisition devices.
7. **Database** – A collection of information organized in such a way that a computer program can quickly select desired pieces of data.
8. **Relational Database Model** – A database system in which any database file can be a component of more than one of the database's tables.
9. **Zoning** – is a partition of the control box of the pattern (i.e. the smallest rectangle containing the pattern); the elements of such partition are used to identify the position in which features of the pattern are detected. It is the process of distinguishing graphic materials from text blocks
10. **Pattern Recognition** – is an important field of computer science concerned with recognizing patterns, particularly visual and sound patterns. It is central to optical character recognition (OCR), voice recognition, and handwriting recognition.
11. **.NET** – is the Microsoft Web services strategy to connect information, people, systems, and devices through software. Integrated across the Microsoft platform, .NET technology provides the ability to quickly build, deploy, manage, and use connected, security-enhanced solutions with Web services.
12. **MICR Toner** – is a special magnetic toner used in printing bank checks on plain paper.