



Aquaforest OCR SDK for .Net

Reference Guide Version 1.10

June 2010

© Copyright 2010 Aquaforest Limited

<http://www.aquaforest.com/>

CONTENTS

1	INTRODUCTION	3
1.1	SDK OVERVIEW	3
1.2	TECHNICAL SUPPORT	3
2	SDK OVERVIEW	4
2.1	SYSTEM REQUIREMENTS	4
2.1.1	<i>Supported Environments</i>	4
2.1.2	<i>.Net Framework</i>	4
2.1.3	<i>Visual C++ Runtime</i>	4
2.2	FOLDERS	4
2.3	A SIMPLE EXAMPLE	4
2.3.1	<i>References</i>	4
2.3.2	<i>Classes</i>	4
2.3.3	<i>Processing Steps</i>	5
2.3.4	<i>C# Example</i>	6
2.3.5	<i>VB.Net Example</i>	7
3	APPLICATION DEVELOPMENT AND DEPLOYMENT	8
3.1	REFERENCES	8
3.2	DEPLOYING C# AND VB.NET APPLICATIONS	8
3.3	DEPLOYING ASP.NET APPLICATIONS	8
3.4	LICENSING	8
3.5	64 BIT SUPPORT	8
4	SAMPLE APPLICATIONS	9
4.1	SIMPLE	9
4.2	CONVERTFOLDEROFFILES	9
4.3	ZONALOCR	9
4.4	ASP.NET	9
5	API REFERENCE	10
5.1	PREPROCESSOR CLASS	10
5.1.1	<i>Constructor</i>	10
5.1.2	<i>Properties</i>	10
5.1.3	<i>Methods</i>	11
5.2	OCR CLASS	12
5.2.1	<i>Constructor</i>	12
5.2.2	<i>Properties</i>	12
5.2.3	<i>Methods</i>	13
5.2.4	<i>Events</i>	15
5.2.5	<i>Enumerations</i>	15
5.3	STATUSUPDATEEVENTARGS CLASS	16
5.3.1	<i>Constructor</i>	16
5.3.2	<i>Properties</i>	16
5.4	WORDS CLASS	17
5.4.1	<i>Constructor</i>	17
5.4.2	<i>Properties</i>	17
5.4.3	<i>Methods</i>	17
5.4.4	<i>WordData Class</i>	18
5.4.5	<i>Properties</i>	18
5.4.6	<i>Error Handling</i>	18
5.4.7	<i>Disposal and Temporary Files folders</i>	18
5.4.8	<i>Multi-threaded applications</i>	18

6	BACKGROUND - SEARCHABLE PDFS	19
6.1	WHAT IS A SEARCHABLE PDF?	19
6.2	INSIDE A SEARCHABLE PDF	19
6.3	OCR ACCURACY	19
6.3.1	<i>Original Image Quality</i>	19
6.3.2	<i>Image DPI and Format</i>	19
6.3.3	<i>Despeckle</i>	19
6.3.4	<i>Deskew</i>	20
6.3.5	<i>Auto-Rotate</i>	20
6.3.6	<i>Graphics Areas</i>	20
6.3.7	<i>Language Settings</i>	20
6.4	HARDWARE AND PERFORMANCE	20
6.4.1	<i>CPU Power</i>	20
6.4.2	<i>Exploiting Multiple CPUs</i>	20
6.4.3	<i>Memory</i>	20
7	ACKNOWLEDGEMENTS	21

1 INTRODUCTION

1.1 SDK Overview

The Aquaforest OCR SDK for .Net incorporates the same high performance OCR engine that is included in the Aquaforest TIFF Junction and Autobahn DX products.

The SDK API allows developers full control over OCR processing to enable customized integration of OCR within .Net applications.

- OCR Bitmap or multi-page TIFF and PDF Files.
- Create Searchable PDF, RTF or Text output files.
- Control pre-processing options such as despeckle, deskew, line removal and autorotate.
- Specify one of the following supported document languages:

- English
- German
- French
- Russian
- Swedish
- Spanish
- Italian
- Russian English
- Ukrainian
- Serbian
- Croatian
- Polish
- Danish
- Portuguese
- Dutch
- Czech
- Roman
- Hungar
- Bulgar
- Slovenian
- Latvian
- Lithuanian
- Estonian
- Turkish

- Enumerate the OCR results, examining the words and characters recognized along with their co-ordinates.
- Process multi-page TIFF and PDF files one page at a time or all in one operation

1.2 Technical Support

Please contact Aquaforest Technical Support with any queries by email support@aquaforest.com. If required, telephone support is also available; please contact Aquaforest using the telephone contact details provided on the company website contact page.

2 SDK OVERVIEW

The SDK is provided as a set of .Net Assemblies, Native DLLs and configuration files designed to allow for straightforward integration into .Net applications.

2.1 System Requirements

2.1.1 Supported Environments

Windows 2003, 2008, Vista, Windows XP.

2.1.2 .Net Framework

.Net Version 3.5

2.1.3 Visual C++ Runtime

The Visual C++ 2008 Redistributable package is required for deployment as well as development.

2.2 Folders

The SDK contains the following folders :

Bin – This contains all the assemblies, DLLs and configurations files.

Docs – SDK Documentation

Samples – Sample C#, VB.Net and ASP.Net samples

2.3 A simple example

The full API reference is in section 5 of this guide, but as a starting point a simple example of a C# and VB.Net console application that creates a searchable PDF from a source TIFF file is shown below.

2.3.1 References

A reference to the Aquaforest.OCR.Api DLL should be added in your application.

If you wish to access the results of the OCR on a word by word basis, for example to obtain word and character results including positional information then you will also need to reference Aquaforest.OCR.Definitions DLL.

2.3.2 Classes

There are two classes used for the OCR :

PreProcessor – This class configures and performs image pre-processing (such as deskewing images) to ensure optimal OCR performance.

Ocr – This is the class that configures and performs the Optical Character Recognition.

Additionally, for accessing the OCR results at an individual word level the following classes are used:

Words – This class contains a collection of words in which is contained all the data available for the words and characters for any given page.

WordData – This class contains a collection of characters that make up the word along with the positional information for each character and the whole word.

StatusUpdateEventArgs – This class is available for each page processed when subscribing to the StatusUpdate event and provides information relating to the processing outcome for the page.

2.3.3 Processing Steps

The following steps are involved in this example

1. Create the `Ocr` and `PreProcessor` objects
2. Specify the location of the OCR bin folder
3. Specify Pre-Processor options
4. Specify OCR Options
5. Read the source file
6. Perform the recognition
7. Save the searchable PDF
8. Delete temporary files (these are by default stored in `%TEMP%` but the location can be specified using `ocr.TempFolder`)

2.3.4 C# Example

```
using System;
using Aquaforest.OCR.Api;

namespace ocr
{
    class Program
    {
        static void Main(string[] args)
        {
            try
            {
                // 1. Create Ocr and Preprocessor Objects
                // and enable console output

                Ocr _ocr = new Ocr();
                PreProcessor _preProcessor = new PreProcessor();
                _ocr.EnableConsoleOutput = true;

                // 2. OCR bin folder Location
                // The bin files can be copied to the application bin
                // folder. Alternatively the System Path and ocr
                // Resource folder can be set as shown below and
                // then just the files in the bin_add folder added
                // to the application bin folder.

                string OCRFiles = @"C:\Aquaforest\OCRSDK\bin";
                System.Environment.SetEnvironmentVariable("PATH",
System.Environment.GetEnvironmentVariable("PATH") + ";" + OCRFiles);

                _ocr.ResourceFolder = OCRFiles;

                // 3. Set PreProcessor Options
                _preProcessor.Deskew = true;
                _preProcessor.Autorotate = false;

                // 4. Set OCR Options
                _ocr.Language = SupportedLanguages.English;
                _ocr.EnablePdfOutput = true;

                // 5. Read Source TIFF File
                _ocr.ReadTIFFSource(
@"C:\Aquaforest\OCRSDK\docs\tiffs\sample.tif");

                // 6. Perform OCR Recognition
                _ocr.Recognize(_preProcessor);

                // 7. Save Output as Searchable PDF
                _ocr.SavePDFOutput(
@"C:\Aquaforest\OCRSDK\docs\tiffs\sample.pdf", true);

                // 8. Clean Up Temporary Files
                _ocr.DeleteTemporaryFiles();
            }
            catch (Exception e)
            {
                Console.WriteLine("Error in OCR Processing : " + e.Message);
            }
        }
    }
}
```

2.3.5 VB.Net Example

```
Module Module1

    Sub Main()
        ' 1. Create Ocr and Preprocessor Objects

        Dim _ocr As New Aquaforest.OCR.Api.Ocr()

        Dim _preProcessor As New Aquaforest.OCR.Api.PreProcessor()

        _ocr.EnableConsoleOutput = True

        ' 2. OCR bin folder Location
        ' The bin files can be copied to the application bin folder.
        ' Alternatively the System Path and ocr Resource folder
        ' can be set as shown below.

        Dim OCRFiles As String
        ' 2. OCR bin folder Location
        ' The bin files can be copied to the application bin
        ' folder. Alternatively the System Path and ocr
        ' Resource folder can be set as shown below and
        ' then just the files in the bin_add folder added
        ' to the application bin folder.

        OCRFiles = "C:\\Aquaforest\\OCRSDK\\bin"
        System.Environment.SetEnvironmentVariable("PATH",
System.Environment.GetEnvironmentVariable("PATH") + ";" + OCRFiles)

        _ocr.ResourceFolder = OCRFiles

        ' 3. Set PreProcessor Options
        _preProcessor.Deskew = True
        _preProcessor.Autorotate = False

        ' 4. Set OCR Options
        _ocr.Language = Aquaforest.OCR.Api.SupportedLanguages.English
        _ocr.EnablePdfOutput = True

        ' 5. Read Source TIFF File
        _ocr.ReadTIFFSource("C:\\Aquaforest\\OCRSDK\\docs\\tiffs\\sample.tif")

        ' 6. Perform OCR Recognition
        _ocr.Recognize(_preProcessor)

        ' 7. Save Output as Searchable PDF
        _ocr.SavePDFOutput("C:\\Aquaforest\\OCRSDK\\docs\\tiffs\\sample.pdf", True)

        '8. Clean Up Temporary Files
        _ocr.DeleteTemporaryFiles()

    End Sub

End Module
```


3 APPLICATION DEVELOPMENT AND DEPLOYMENT

3.1 References

To use the API a reference to Aquaforest.Ocr.Api must be included in your application. If you wish to enumerate the OCR results rather than simply generate PDF, RTF or TXT outputs then you will also need to add a reference to Aquaforest.Ocr.Definitions.

3.2 Deploying C# and VB.Net Applications

Any deployment method should ensure that the target system meets the requirements (see section 2.1) and should include the full contents of the SDK *bin* folder.

There are two approaches to handling the contents of the SDK *bin* folder that can be used when building and deploying C# and VB.Net applications.

Approach 1

The entire contents of the SDK *bin* folder can be copied to the application bin folder.

Approach 2

The contents of the SDK *bin_add* folder can be copied to the application bin folder and the ocr resource folder and path of the full bin folder must be specified as shown in the sample code in section 2 above. In certain cases, e.g. deploying as a service, this approach is not suitable.

3.3 Deploying ASP.Net Applications

The same two approaches that work for C# and VB.Net can also be employed for ASP.Net applications. Note that with trial licenses a pop-up dialog box appears on the server.

3.4 Licensing

Production system deployment requires that a license string is defined in the code. The license string defines the number of concurrent OCR processes that can be run.

For example :

```
ocr.License =  
"MT0xMjM0NTY7BLk4uT3RoZXOzM9NDs0PVRydWEYzMDRFOEQxMzg0QkQ5ODREQTk3RQ" ;
```

If the string is not specified the SDK will run in evaluation mode. In evaluation mode :

- *A trial “pop-up” will appear for each document processed*
- *Generated searchable PDFs will include indelible watermarks*
- *Only 3 pages are generated for text or RTF files.*

3.5 64 bit support

Whilst applications built using the SDK will run on 64 bit systems, they must currently be built as 32 bit.

4 SAMPLE APPLICATIONS

The samples folder includes a number of sample applications in C#, VB.Net and ASP.Net. The project files are for Visual Studio 2008. The sample applications are as described below.

4.1 Simple

This includes VB and C# version of the simple application listed in section 2.

4.2 ConvertFolderOfFiles

This forms-based application demonstrates converting a folder of TIFF or Bitmap files. Both VB.Net and C# versions are included.

4.3 ZonalOCR

This forms-based application demonstrates converting a folder of TIFF or Bitmap files. Both VB.Net and C# versions are included.

4.4 ASP.Net

This demonstrates a simple web-based application that allows uploading an image file for OCR and conversion to text, RTF or searchable PDF.

5 API REFERENCE

To use the API a reference to `Aquaforest.Ocr.Api` must be included in your application. If you wish to enumerate the OCR results rather than simply generate PDF, RTF or TXT outputs then you will also need to add a reference to `Aquaforest.Ocr.Definitions`.

5.1 PreProcessor Class

A `PreProcessor` object, which must be created and passed to the `Ocr` object, controls all of the pre-processing that can be performed on the input image in order to improve the quality of the output. Instantiation of the `PreProcessor` object will initialise a default set of pre-processing options which result in minimal image manipulation. For a full description of the pre-processing options available and appropriate values see section 5.1.2 Properties below.

5.1.1 Constructor

```
PreProcessor preProcessor = new PreProcessor();
```

5.1.2 Properties

Property	Description
bool Autorotate	Auto-rotate the image – this will ensure all text oriented normally. The default value is false (disabled). Note: When using a PDF source Autorotation will be disabled on any pages already containing text.
int Binarize	This value should generally only be used under guidance from technical support. It can control the way that color images are processed and force binarization with a particular threshold. A value of 200 has been shown to generally give good results in testing, but this should be confirmed with “typical” customer documents.
int BlankPageThreshold	Use this to set the minimum number of "On Pixels" that must be present in the image for a page not to be considered blank. A value of -1 will turn off blank page detection. A value of 100 produced reasonable blank page detection in testing, but the validity of this should be confirmed using “typical” source documents.
int BoxSize	This option is ideal for forms where sometimes boxes around text can cause an area to be identified as graphics. This option removes boxes from the temporary copy of the imaged used by the OCR engine. It does not remove boxes from the final image. Technically, this option removes connected elements with a minimum area (in pixels and defined by this property). This option is currently only applied for bitonal images.
int Despeckle	Despeckle the image – The method removes all disconnected elements within the image that have height or width in pixels less than the specified figure. The maximum value is 9 and the default value is 0.
bool Deskew	Deskew (straighten) the image. The default value is false (disabled).
string Morph	Image Morphology. This option should generally only be used under guidance from technical support.

PreProcessor Class Properties - Continued

Property	Description
bool NoPictures	By default, if an area of the document is identified as a graphic area then no OCR processing is run on that area. However, certain documents may include areas or boxes that are identified as “graphic” or “picture” areas but that actually do contain useful text. Setting NoPictures to false will cause it to ignore areas identified as pictures whilst setting it to true will force OCR of areas identified as pictures.
bool RemoveLines	When set to true this will enable the removal of lines. This feature is particularly useful where pages contain tables and underlining which can prevent the OCR engine from recognising characters. The lines are removed only from the image used in OCR and not from the image used in the final PDF if PDF creation is enabled.

5.1.3 Methods

Method	Description
ConfigurePDFStamp(string prefix, string suffix, Nullable<int> start, Nullable<int> digits, PagePositionEnum position, StampType stampType)	<p>Using this method stamps can be configured to be added to each page of the PDF output. The stamps contain one or more of the following:</p> <ul style="list-style-type: none"> • Prefix – a string to be added to the beginning of the stamp, before the number section. • Start – the value that the number portion of the stamp should start at. The number portion will be incremented by 1 each page. • Digits – a value indicating the minimum length that the number portion of the stamp should be displayed as. Preceding 0’s will be used to pad any numbers less than this whilst numbers greater than this will be displayed in full. • Suffix - a string to be added to the end of the stamp, after the number section. <p>Thus a stamp with Prefix = “Beginning”, Start = “1”, Digits = “4” and Suffix = “End” would produce the text “Beginning0001End” on the first page. Any one of these can be set to null resulting in the exclusion of that part from the final text.</p> <p>Additionally the stamp can be added either as visible searchable text or as an image and can be positioned in one of the following:</p> <ul style="list-style-type: none"> • Top Left • Top Centre • Top Right • Centre Left • Centre • Centre Right • Bottom Left • Bottom Centre • Bottom Right

5.2 OCR Class

The OCR object is used to control OCR processing, obtain status updates during processing and retrieve the resulting output from this processing upon completion.

5.2.1 Constructor

```
Ocr ocr = new Ocr();
```

5.2.2 Properties

Property	Description
String ResourceFolder	This property can optionally be used to set the location of the resources folder when the resources are not located in the same folder as the assembly using the API.
SupportedLanguages Language	Sets the language to be used for the OCR processing. This takes a value from the enumeration SupportedLanguages which is defined in the API. Default language is English.
bool EnablePDFOutput	Enables or disables the production of Portable Document Format output. Default value is false (disabled).
bool EnableTextOutput	Enables or disables the production of simple text final output. Default value is true (enabled).
bool EnableRTFOutput	Enables or disables the production of Rich Text Format output. Default value is false (disabled).
int StartPage	Sets the first page of the source file that the OCR process will be begin from (for a multipage source). Throws an <code>ArgumentOutOfRangeException</code> if a source file has not been set already (by using the <code>ReadBMPSource</code> or <code>ReadTIFFSource</code> method prior to setting this property) or if the page is greater than the number of pages in the source. By default the whole of the document will be processed.
int EndPage	Sets the last page of the source file that the OCR process will be run to (for a multipage source). Throws an <code>ArgumentOutOfRangeException</code> if a source file has not been set already (by using the <code>ReadBMPSource</code> or <code>ReadTIFFSource</code> method prior to setting this property) or if the page is greater than the number of pages in the source. By default the whole of the document will be processed.
int CurrentPage	Returns the current page for which the OCR has been performed. This is useful only when using <code>Recognize()</code> in another thread.
bool HandleExceptionsInternally	When set to true the <code>Ocr</code> object will catch any exceptions for method calls and simply return false from the method. The exceptions caught are stored in the <code>LastException</code> property overwriting any previous value.
Exception LastException	Stores last exception caught by the <code>Ocr</code> object.
bool EnableConsoleOutput	If enabled then progress messages will be sent to the console. Default is false.

OCR Class Properties - Continued

Property	Description
string TempFolder	Specifies a temporary folder for storing bitmap images and intermediate output during OCR processing. If this is not specified, the first of the following environment variables that is defined will be used : TMP, TMPDIR, TEMP.
String License	Specifies the license key
EnableConsoleOutput	If set to True, progress messages will be written to the console output. Default false.
EnableDebugOutput	If set to True, debug messages will be written to the console output. Default false.
RemoveExistingPDFText	RemoveExistingPDFText if set to true will result in the removal of any existing text from the output PDF*. *Note: when PDF output is generated from a PDF source it is a copy of the PDF that is manipulated rather than generating a new one. This approach offers several advantages such as potential size savings and performance enhancements.

5.2.3 Methods

Method	Description
void ReadBMPSource(string fileName)	Checks for the existence of the source file and sets up the OCR engine for handling the bitmap image.
void ReadPDFSource(string fileName)	Checks for the existence of the source file and sets up the OCR engine for handling the PDF.
void ReadPDFSource(string filename, string password)	Checks for the existence of the source file and sets up the OCR engine for handling the secure PDF for which the password is provided. If PDF output is generated from this the output will have no security settings defined.
void ReadTIFFSource(string fileName)	Checks for the existence of the source file and sets up the OCR engine for handling the TIFF image.
void Recognize(PreProcessor preProcessor)	Performs any pre-processing defined in the PreProcessor object and then carries out OCR processing on the pre-processed image.
Words ReadPageWords(int pagenumber)	Returns an instance of the Words class for the specified page.
Words ReadPageWords(int pagenumber, Rectangle region)	Returns an instance of the Words class for the specified page where the words are fully enclosed in the bounds of the region specified.
string ReadPageString(int pagenumber)	Returns a string containing the words from the

	specified page.
--	-----------------

OCR Class Methods - Continued

Method	Description
string ReadPageString(int pageNumber, Rectangle region)	Returns a string containing the words for the specified page where the words are fully enclosed in the bounds of the region specified.
bool SavePDFOutput(string fileName, bool overwriteExisting)	Saves the output to a PDF file with the name specified. If any text was extracted then this will be searchable in the PDF.
bool SaveRTFOutput(string fileName, bool overwriteExisting)	Saves the output to a RTF file with the name provided.
bool SaveTextOutput(string fileName, bool overwriteExisting)	Saves the text extracted to a simple text file with the name provided.
Image GetPageImage(int pageNumber)	Returns a System.Drawing.Image containing the processed image.
void DeleteTemporaryFiles()	Removes temporary files created during the OCR processing from the system. Note, do not call this before you have completely finished processing a file.
void Abort()	Stops processing of an ongoing call to Recognize.

5.2.4 Events

Event	Description
void PageCompleted(int pageNumber, bool ocrSuccess, bool imageProcessingSuccess, bool blankPage)	This event is raised when processing of a page is complete. The variable pageNumber indicates which page has just completed, the variable ocrSuccess indicates whether the OCR for that page was successful or not (note, a successful OCR does not necessarily indicate that text was found on a page), imageProcessingSuccess indicates whether the pre-process was successful for that page, blankPage indicates whether the page detected as blank.
void StatusUpdate (object sender, StatusUpdateEventArgs statusUpdateEventArgs)	This event is raised when processing of a page is complete. The StatusUpdateEventArgs object provides access to information relating to the status of the page processed.

5.2.5 Enumerations

Enumeration	Description
SupportedLanguages	This enumeration includes all of the languages currently supported by the API.

5.3 StatusUpdateEventArgs Class

This class contains information relating to the conversion status of a page.

5.3.1 Constructor

An instance of this class is obtained for each page processed when subscribing to the event StatusUpdate.

5.3.2 Properties

Property	Description
int PageNumber	This property returns page for which the object relates to.
int Rotation	A value from 0 to 3 which indicates the rotation used for the output in terms of the number of 90° steps away from the orientation in which the input page was provided. If AutoRotation is set to false this will always be 0.
double ConfidenceScore	Generally a value of 1 or greater would indicate that reasonable OCR of a page, but this should be confirmed using “typical” source files.
bool TextAvailable	This property indicates whether text was extracted for the page.
bool ImageAvailable	This property indicates whether an image (after all appropriate pre-processing) was successfully extracted.
bool BlankPage	This property indicates whether the page was detected as blank.

5.4 Words Class

This class contains a collection of WordData objects which are available on a page by page basis.

5.4.1 Constructor

An instance of this class is obtained by calling the ReadPageWords method on the Ocr object, passing the page for which the words are required.

5.4.2 Properties

Property	Description
int Count	This property returns the number of WordData objects in the collection.
int Height	This property returns the height of the current word.
int Width	This property returns the width of the current page word.

5.4.3 Methods

Method	Description
WordData GetFirst()	Returns the first WordData object in the collection and sets the index to this item.
WordData GetNext()	Returns the next WordData object in the collection and sets the index to this item.
int GetHeight(int index)	Returns the word height from the WordData object stored at the specified index in the collection.
int GetWidth(int index)	Returns the word width from the WordData object stored at the specified index in the collection.

5.4.4 WordData Class

This class contains the individual characters along with the positional information relating to each character in the word and to the word as a whole.

5.4.5 Properties

Property	Description
float AverageCharacterHeight	This property returns the average height of all the characters in the word.
float AverageCharacterWidth	This property returns the average width of all the characters in the word.
int Bottom	This property returns the bottom of the word.
int CharacterList	This property returns a list of CharacterData objects for the word.
int Height	This property returns the height of the word.
int Left	This property returns the left edge of the word.
int Top	This property returns the Top of the word.
int Width	This property returns the width of the word.
string Word	This property returns the word as a string.

5.4.6 Error Handling

There are two options regarding error handling using the API.

1. Using the default settings various exceptions can be thrown by the Ocr object so these should be trapped within the calling code.
2. Alternatively HandleExceptionsInternally can be set to true with the result that method calls will return false on error but throw no exceptions. The calling code can obtain the last exception from the LastException property if details of the failure are required.

5.4.7 Disposal and Temporary Files folders

During the OCR processing various temporary files are generated and used at different stages. These temporary files can be removed by calling DeleteTemporaryFiles. However, such a call should not be made until all processing (both within the Ocr object and calling code) on a file is complete as these files are required when calling SaveRTFOutput, SavePDFOutput, SaveTextOutput, GetPageImage and ReadPageWords. When the Ocr object is disposed of the temporary files are automatically removed.

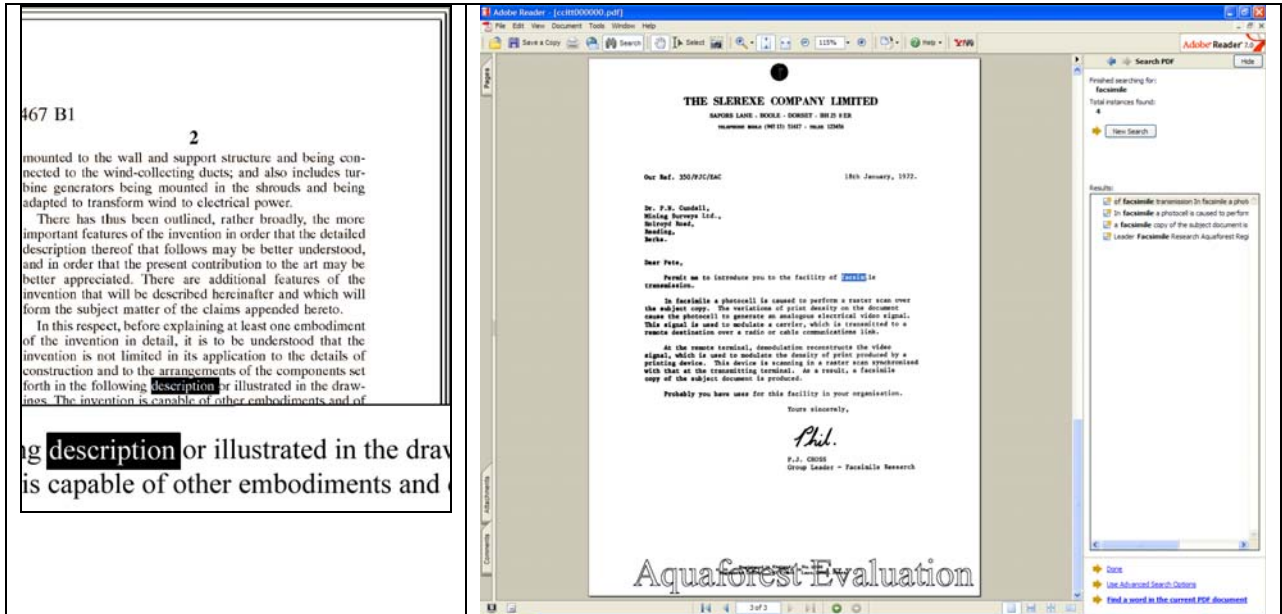
5.4.8 Multi-threaded applications

Temporary files created and used throughout the OCR processing are named according to the page number, therefore if Ocr objects are instantiated in multiple threads then a different temporary folder must be set for each folder. If this is not done then un-expected behaviour will result.

6 BACKGROUND - SEARCHABLE PDFS

6.1 What is a Searchable PDF?

A searchable PDF file is a PDF file that includes text that can be searched upon using the standard Adobe Reader “search” functionality. In addition, the text can be selected and copied from the PDF. Generally, PDF files created from Microsoft Office Word and other documents are by their nature searchable as the source document contains text which is replicated in the PDF, but when creating a PDF from a scanned document and an OCR process needs to be applied to recognize the characters within the image.



6.2 Inside a Searchable PDF

In the context of Document Imaging, a searchable PDF will typically contain both the original scanned image plus a separate text layer produced from an OCR process. The text layer is defined in the PDF file as invisible, but can still be selected and searched upon. PDF files are able to store images using most of the native compression schemes used in TIFF files, so for example Group 4 TIFF files do not usually require any format conversion.

6.3 OCR Accuracy

A number of factors affect the accuracy of the text produced by the OCR process – 100% accuracy is certainly possible under good conditions but each of the following issues, and OCR processing options will have an impact.

6.3.1 Original Image Quality

Although some pre-processing options such as despeckle and deskew can help in some cases, the visual quality of the original scan is of paramount importance.

6.3.2 Image DPI and Format

The image resolution should be at least 150 DPI for OCR processing, and preferably 300 DPI for optimal results, although for good quality scans 200 DPI is often sufficient. Non-lossy formats (TIFF Group 4, LZW etc) are preferred over lossy formats such as JPEG.

6.3.3 Despeckle

This pre-processing option removes isolated “dots” within the image which can cause recognition problems, and makes the result image “cleaner”.

6.3.4 Deskew

This option can improve OCR results by straightening crooked pages.

6.3.5 Auto-Rotate

OCR processing usually recognizes text written top-to-bottom, left-to-right, so pages that are orientated any other way (usually landscape pages) need to be re-oriented to enable recognition.

6.3.6 Graphics Areas

There are two options that can be used to control how the OCR engine processes parts of the document image that appear to be graphics areas.

To ensure that the OCR engine can be forced to process such areas there are two options :

“Treat all Graphics Areas as Text”. This option will ensure the entire document is processed as text.

“Remove Box Lines in OCR Processing”. This option is ideal for forms where sometimes boxes around text can cause an area to be identified as graphics. This option removes boxes from the temporary copy of the imaged used by the OCR engine. It does not remove boxes from the final image. Technically, this option removes connected elements with a minimum area (by default 100 pixels).

6.3.7 Language Settings

The language setting determines the set of characters that will be recognized, and the dictionary that will be used as a guide.

6.4 Hardware and Performance

6.4.1 CPU Power

The OCR process is highly CPU intensive and will benefit from being given as much CPU power as possible. As a guide about 2,000 pages per hour can be processed on a 3.0 GHz processor core, although this will vary according to the source document and OCR options chosen.

6.4.2 Exploiting Multiple CPUs

To take advantage of multiple cores, multiple OCR instances should be run in parallel.

6.4.3 Memory

Memory can be a limiting factor when creating the final PDF, in the case of very large documents. A rule of thumb would be to have 1GB – 1.5 GB of memory per processor core.

7 ACKNOWLEDGEMENTS

This product makes use of a number of Open Source components which are included in binary form. The appropriate acknowledgements and copyright notices are given below.

LEPTONICA

Copyright (C) 2001 Leptonica. All rights reserved.

LIBJPEG

This software is based in part on the work of the Independent JPEG Group.

ZLIB

(C) 1995-2004 Jean-loup Gailly and Mark Adler.

ITEXT 4.1.6

Copyright (C) 1999-2009 by Bruno Lowagie and Paulo Soares et al. All Rights Reserved. Binaries distributed under the Mozilla Public License.

CUNEIFORM

Copyright (c) 1993-2008, Cognitive Technologies. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met

Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer. Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution. Neither the name of the Cognitive Technologies nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE

LIBTIFF

Copyright (c) 1988-1997 Sam Leffler. Copyright (c) 1991-1997 Silicon Graphics, Inc.

Permission to use, copy, modify, distribute, and sell this software and its documentation for any purpose is hereby granted without fee, provided that (i) the above copyright notices and this permission notice appear in all copies of the software and related documentation, and (ii) the names of Sam Leffler and Silicon Graphics may not be used in any advertising or publicity relating to the software without the specific, prior written permission of Sam Leffler and Silicon Graphics.

THE SOFTWARE IS PROVIDED "AS-IS" AND WITHOUT WARRANTY OF ANY KIND, EXPRESS, IMPLIED OR OTHERWISE, INCLUDING WITHOUT LIMITATION, ANY WARRANTY OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. IN NO EVENT SHALL SAM LEFFLER OR SILICON GRAPHICS BE LIABLE FOR ANY SPECIAL, INCIDENTAL, INDIRECT OR CONSEQUENTIAL DAMAGES OF ANY KIND, OR ANY DAMAGES WHATSOEVER RESULTING FROM LOSS OF USE, DATA OR PROFITS, WHETHER OR NOT ADVISED OF THE POSSIBILITY OF DAMAGE, AND ON ANY THEORY OF LIABILITY, ARISING OUT OF OR IN CONNECTION WITH THE USE OR PERFORMANCE OF THIS SOFTWARE.

FREEIMAGE

This software uses the FreeImage open source image library. See <http://freeimage.sourceforge.net> for details. FreeImage is used under the (FIPL), version 1.0.