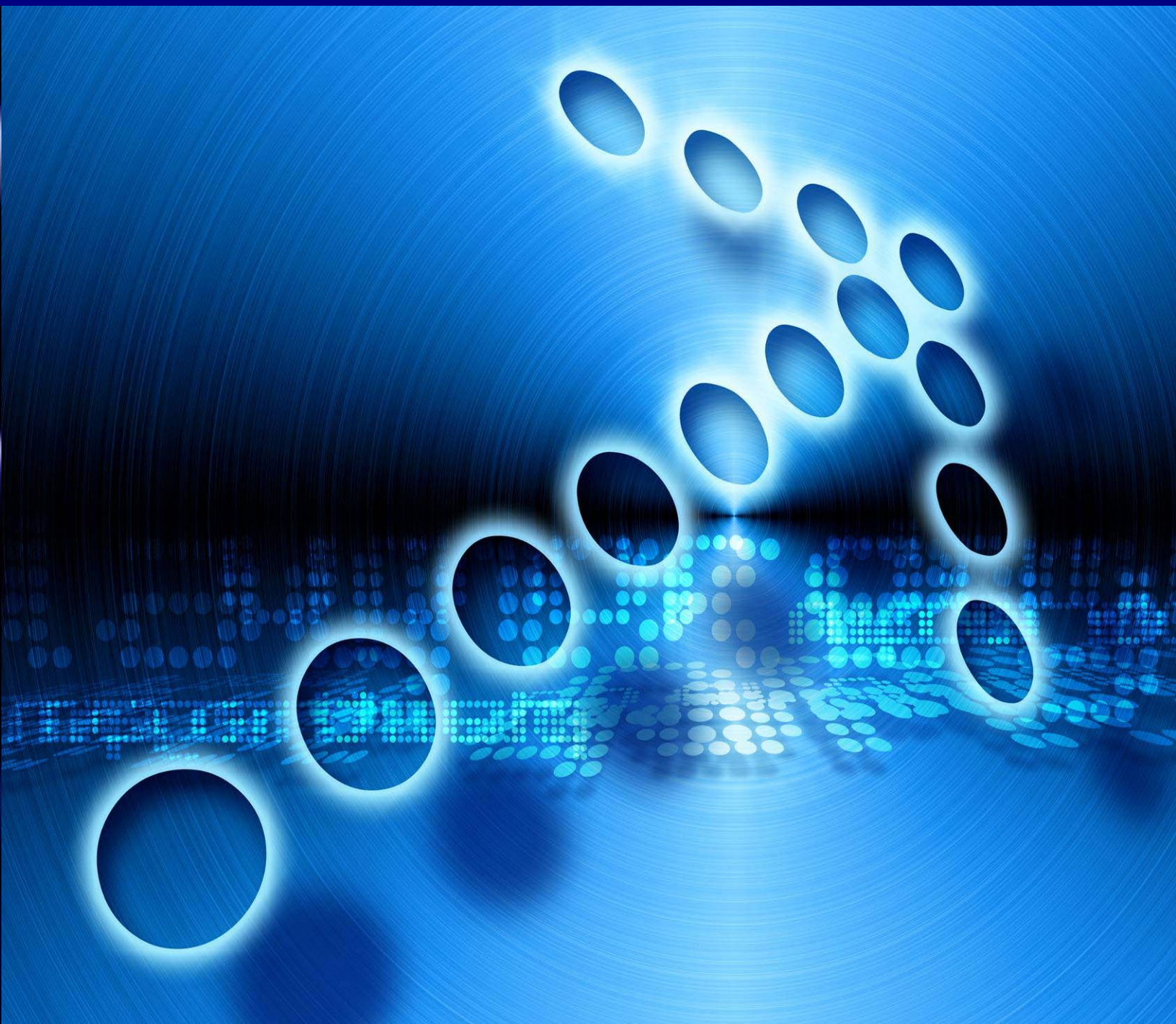


# Creating Searchable PDFs From Scanned Documents

**Aquaforest Limited**

[www.aquaforest.com](http://www.aquaforest.com)



# Searchable PDF Explained

This brief document aims to provide guidance for the creation of searchable PDF files from scanned documents, whether standard TIFF Files or Image-Only PDF files.

## Types of PDF File

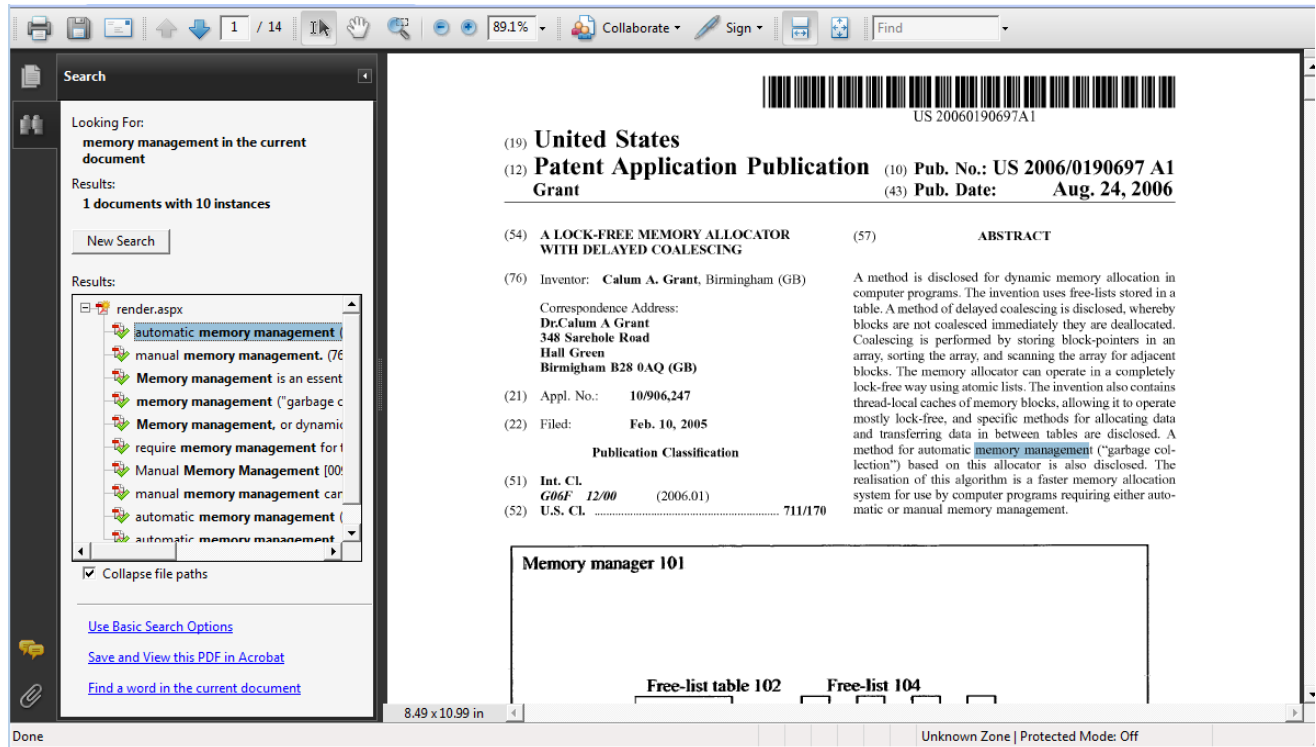
PDF Type	Description
Normal	This is the most common type of PDF and is most typically created from a document such as Microsoft Word. It contains the full text of the page with appropriate coding to define fonts, sizes, etc. and will provide a faithful print of the original.
Image Only	This is a PDF that has been created from one or more images – most commonly as a result of scanning a document either directly to PDF or by converting a scanned TIFF image to PDF. These files do not contain any searchable text and most often comprise a set of Group4 or JBIG2 images in a PDF “wrapper”.
Searchable	A “Searchable” PDF is an “Image-Only” PDF that additionally contains a hidden layer of text generated by an OCR engine. This enables the file to be searched in the same fashion as a “Normal” PDF. Text can be copied and pasted.

## 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.

# 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.



Searching a Searchable PDF with Adobe Reader

## 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

## 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 compression

### **Despeckle**

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

### **Deskew**

This option can improve OCR results by straightening crooked pages.

### **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.

### **Language Settings**

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



# Hardware and Performance

## 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 1,000 pages per hour can be processed on a 2.5GHz processor, although this will vary according to the source document and OCR options chosen.

## Exploiting Multiple CPUs

To take advantage of multiple CPUs, multiple conversion jobs should be run concurrently. This can most conveniently be done by using the Job Management facilities of Autobahn DX [2]

## 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.

# References

[1] TIFF Junction

[http://www.aquaforest.com/en/merge\\_tiff\\_junction.asp](http://www.aquaforest.com/en/merge_tiff_junction.asp)

[2] Autobahn DX

<http://www.aquaforest.com/en/autobahn.asp>