# scanned PDF to text in C# and ByteScout PDF Extractor SDK

## How to use ByteScout PDF Extractor SDK for scanned PDF to text in C#

We've created and updating regularly our sample code library so you may quickly learn scanned PDF to text and the step-by-step process in C#. ByteScout PDF Extractor SDK was made to help with scanned PDF to text in C#. ByteScout PDF Extractor SDK is the Software Development Kit (SDK) that is designed to help developers with data extraction from unstructured documents like pdf, tiff, scans, images, scanned and electronic forms. The library is powered by OCR, computer vision and AI to provide unique functionality like table detection, automatic table structure extraction, data restoration, data restructuring and reconstruction. Supports PDF, TIFF, PNG, JPG images as input and can output CSV, XML, JSON formatted data. Includes full set of utilities like pdf splitter, pdf merger, searchable pdf maker.

The SDK samples like this one below explain how to quickly make your application do scanned PDF to text in C# with the help of ByteScout PDF Extractor SDK. In order to implement this functionality, you should copy and paste code below into your app using code editor. Then compile and run your application. Test C# sample code examples whether they respond your needs and requirements for the project.

Trial version can be downloaded from our website. Source code samples for C# and documentation are included.

FOR MORE INFORMATION AND FREE TRIAL:

Download Free Trial SDK (on-premise version)

Read more about ByteScout PDF Extractor SDK

Explore API Documentation

Get Free Training for ByteScout PDF Extractor SDK

Get Free API key for Web API

visit www.ByteScout.com

## Source Code Files:

Program.cs

```csharp
using System.Diagnostics;
using Bytescout.PDFExtractor;

// This example demonstrates the use of Optical Character Recognition (OCR) to extract
// from scanned PDF documents and raster images.

// To make OCR work you should add the following references to your project:
// 'Bytescout.PDFExtractor.dll', 'Bytescout.PDFExtractor.OCRExtension.dll'.

namespace ScannedPdfToText
{
    class Program
    {
        static void Main(string[] args)
        {
            // Create Bytescout.PDFExtractor.TextExtractor instance
            TextExtractor extractor = new TextExtractor();
            extractor.RegistrationName = "demo";
            extractor.RegistrationKey = "demo";

            // Load sample PDF document
            extractor.LoadDocumentFromFile("sample_ocr.pdf");

            // Enable Optical Character Recognition (OCR)
            // in .Auto mode (SDK automatically checks if needs to use OCR or not)
            extractor.OCRMode = OCRMode.Auto;

            // Set the location of OCR language data files
            extractor.OCRLanguageDataFolder = @"c:\Program Files\Bytescout PDF Extracto

            // Set OCR language
            extractor.OCRLanguage = "eng"; // "eng" for english, "deu" for German, "fr
            // Find more language files at https://github.com/bytescout/ocrdata

            // Set PDF document rendering resolution
            extractor.OCRResolution = 300;


            // You can also apply various preprocessing filters
            // to improve the recognition on low-quality scans.

            // Automatically deskew skewed scans
            //extractor.OCRImagePreprocessingFilters.AddDeskew();

            // Remove vertical or horizontal lines (sometimes helps to avoid OCR engine
            //extractor.OCRImagePreprocessingFilters.AddVerticalLinesRemover();
            //extractor.OCRImagePreprocessingFilters.AddHorizontalLinesRemover();

            // Repair broken letters
            //extractor.OCRImagePreprocessingFilters.AddDilate();

            // Remove noise
            //extractor.OCRImagePreprocessingFilters.AddMedian();

            // Apply Gamma Correction
```

```csharp
        //extractor.OCRImagePreprocessingFilters.AddGammaCorrection();

        // Add Contrast
        //extractor.OCRImagePreprocessingFilters.AddContrast(20);


        // (!) You can use new OCRAnalyser class to find an optimal set of image pr
        // filters for your specific document.
        // See "OCR Analyser" example.


        // Save extracted text to file
        extractor.SaveTextToFile("output.txt");

        // Cleanup
        extractor.Dispose();

        // Open result document in default associated application (for demo purpos
        ProcessStartInfo processStartInfo = new ProcessStartInfo("output.txt");
        processStartInfo.UseShellExecute = true;
        Process.Start(processStartInfo);
    }
  }
}
```

VIDEO

https://www.youtube.com/watch?v=s28W3_KMraU

ON-PREMISE OFFLINE SDK

60 Day Free Trial or Visit ByteScout PDF Extractor SDK Home Page
Explore ByteScout PDF Extractor SDK Documentation
Explore Samples
Sign Up for ByteScout PDF Extractor SDK Online Training

ON-DEMAND REST WEB API

Get Your API Key
Explore Web API Docs
Explore Web API Samples

[visit www.ByteScout.com](www.ByteScout.com)

[visit www.PDF.co](www.PDF.co)