



PyMuPDF Documentation

Release 1.11.0

Jorj X. McKie

Jul 30, 2017

CONTENTS

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 1.1 | Note on the Name <code>fitz</code> | 2 |
| 1.2 | License | 2 |
| 1.3 | Covered Version | 2 |
| 2 | Installation | 3 |
| 2.1 | Option 1: Install from Sources | 3 |
| 2.1.1 | Step 1: Download PyMuPDF | 3 |
| 2.1.2 | Step 2: Download and Generate MuPDF | 3 |
| 2.1.3 | Step 3: Build / Setup PyMuPDF | 5 |
| 2.1.4 | Using UPX | 5 |
| 2.2 | Option 2: Install from Binaries | 5 |
| 2.2.1 | Step 1: Download Optional Material | 5 |
| 2.2.2 | Step 2: Install PyMuPDF | 5 |
| 2.2.3 | MD5 Checksums | 6 |
| 2.2.4 | Targeting Parallel Python Installations | 6 |
| 3 | Tutorial | 7 |
| 3.1 | Importing the Bindings | 7 |
| 3.2 | Opening a Document | 7 |
| 3.3 | Some Document Methods and Attributes | 7 |
| 3.4 | Accessing Meta Data | 8 |
| 3.5 | Working with Outlines | 8 |
| 3.6 | Working with Pages | 8 |
| 3.6.1 | Inspecting the Links of a Page | 8 |
| 3.6.2 | Rendering a Page | 9 |
| 3.6.3 | Saving the Page Image in a File | 9 |
| 3.6.4 | Displaying the Image in Dialog Managers | 9 |
| 3.6.5 | Extracting Text | 9 |
| 3.6.6 | Searching Text | 10 |
| 3.7 | PDF Maintenance | 10 |
| 3.7.1 | Modifying, Creating, Re-arranging and Deleting Pages | 10 |
| 3.7.2 | Joining and Splitting PDF Documents | 11 |
| 3.7.3 | Saving | 11 |
| 3.8 | Closing | 11 |
| 3.9 | Example: Dynamically Cleaning up Corrupt PDF Documents | 12 |
| 4 | Classes | 13 |
| 4.1 | Annot | 13 |
| 4.1.1 | Example | 17 |
| 4.2 | Colorspace | 17 |
| 4.3 | Document | 18 |
| 4.3.1 | Remarks on <code>select()</code> | 29 |
| 4.3.2 | <code>select()</code> Examples | 29 |

| | | |
|----------|---|-----------|
| 4.3.3 | <code>setMetadata()</code> Example | 30 |
| 4.3.4 | <code>setToC()</code> Example | 31 |
| 4.3.5 | <code>insertPDF()</code> Examples | 31 |
| 4.3.6 | Other Examples | 31 |
| 4.4 | Identity | 32 |
| 4.5 | IRect | 32 |
| 4.5.1 | Remark | 35 |
| 4.5.2 | IRect Algebra | 35 |
| 4.5.3 | Examples | 36 |
| 4.6 | Link | 36 |
| 4.7 | <code>linkDest</code> | 37 |
| 4.8 | Matrix | 39 |
| 4.8.1 | Remarks 1 | 42 |
| 4.8.2 | Remarks 2 | 42 |
| 4.8.3 | Matrix Algebra | 42 |
| 4.8.4 | Examples | 43 |
| 4.8.5 | Shifting | 43 |
| 4.8.6 | Flipping | 44 |
| 4.8.7 | Shearing | 45 |
| 4.8.8 | Rotating | 46 |
| 4.9 | Outline | 47 |
| 4.10 | Page | 48 |
| 4.10.1 | Common Parameters | 56 |
| 4.10.2 | Description of <code>getLinks()</code> Entries | 58 |
| 4.10.3 | Notes on Supporting Links | 58 |
| 4.10.4 | Homologous Methods of Document and Page | 59 |
| 4.11 | Pixmap | 59 |
| 4.11.1 | Supported Pixmap Construction Image Types | 64 |
| 4.11.2 | Details on Saving Images with <code>writeImage()</code> | 64 |
| 4.11.3 | Pixmap Example Code Snippets | 64 |
| 4.12 | Point | 67 |
| 4.12.1 | Remark | 68 |
| 4.12.2 | Point Algebra | 68 |
| 4.12.3 | Examples | 68 |
| 4.13 | Rect | 69 |
| 4.13.1 | Remark | 73 |
| 4.13.2 | Rect Algebra | 73 |
| 4.13.3 | Examples | 73 |
| 5 | Low Level Functions and Classes | 75 |
| 5.1 | Functions | 75 |
| 5.2 | Device | 80 |
| 5.3 | DisplayList | 80 |
| 5.4 | TextPage | 81 |
| 5.5 | TextSheet | 82 |
| 5.6 | Working together: Device, DisplayList, TextPage and TextSheet | 82 |
| 5.6.1 | Generate Pixmap | 82 |
| 5.6.2 | Perform Text Search | 82 |
| 5.6.3 | Extract Text | 83 |
| 6 | Constants and Enumerations | 85 |
| 6.1 | Constants | 85 |
| 6.2 | Text Alignment | 86 |
| 6.3 | Enumerations | 86 |
| 6.4 | Link Destination Flags | 86 |
| 6.5 | Annotation Types | 87 |
| 6.6 | Annotation Flags | 88 |
| 6.7 | Annotation Line End Styles | 89 |

| | | |
|-----------|---|------------|
| 7 | Color Database | 91 |
| 7.1 | Function <code>getColor()</code> | 91 |
| 7.2 | Printing the Color Database | 91 |
| 8 | Appendix 1: Performance | 93 |
| 8.1 | Part 1: Parsing | 93 |
| 8.2 | Part 2: Text Extraction | 97 |
| 8.3 | Part 3: Image Rendering | 98 |
| 9 | Appendix 2: Details on Text Extraction | 101 |
| 9.1 | General structure of a <code>TextPage</code> | 101 |
| 9.2 | Output of <code>getText(output="text")</code> | 101 |
| 9.3 | Output of <code>getText(output="html")</code> | 102 |
| 9.4 | Output of <code>getText(output="json")</code> | 102 |
| 9.5 | Output of <code>getText(output="xml")</code> | 103 |
| 9.6 | Performance | 103 |
| 10 | Appendix 3: Considerations on Embedded Files | 105 |
| 10.1 | General | 105 |
| 10.2 | MuPDF Support | 105 |
| 10.3 | PyMuPDF Support | 105 |
| 11 | Appendix 4: Assorted Technical Information | 107 |
| 11.1 | PDF Base 14 Fonts | 107 |
| 11.2 | Adobe PDF Reference 1.7 | 107 |
| 12 | Change Logs | 109 |
| 12.1 | Changes in Version 1.11.0 | 109 |
| 12.2 | Changes in Version 1.10.0 | 110 |
| 12.2.1 | MuPDF v1.10 Impact | 110 |
| 12.2.2 | Other Changes compared to Version 1.9.3 | 110 |
| 12.3 | Changes in Version 1.9.3 | 111 |
| 12.4 | Changes in Version 1.9.2 | 111 |
| 12.5 | Changes in Version 1.9.1 | 112 |
| 13 | Error Messages | 115 |

INTRODUCTION



PyMuPDF is a Python binding for [MuPDF](#) - “a lightweight PDF and XPS viewer”.

MuPDF can access files in PDF, XPS, OpenXPS, CBZ (comic book archive), FB2 and EPUB (e-book) formats.

These are files with extensions `*.pdf`, `*.xps`, `*.oxps`, `*.cbz`, `*.fb2` or `*.epub` (so in essence, with this binding you can develop **e-book viewers in Python** ...)

PyMuPDF provides access to many important functions of MuPDF from within a Python environment, and we are continuously seeking to expand this function set.

MuPDF stands out among all similar products for its top rendering capability and unsurpassed processing speed. At the same time, its “lightweight” makes it an excellent choice for platforms where resources are typically limited, like smartphones.

Check this out yourself and compare the various free PDF-viewers. In terms of speed and rendering quality [SumatraPDF](#) ranges at the top (apart from MuPDF’s own standalone viewer) - since it has changed its library basis to MuPDF!

While PyMuPDF has been available since several years for an earlier version of MuPDF (v1.2, called **fitz-python** then), it was until only mid May 2015, that its creator and a few co-workers decided to elevate it to support current releases of MuPDF (first v1.7a, up to latest 1.11.0 in May 2017).

PyMuPDF runs and has been tested on Mac, Linux, Windows XP SP2 and up, Python 2.7 through Python 3.6 (only up to 3.4 on Windows XP), 32bit and 64bit versions. Other platforms should work too, if MuPDF and Python support them.

PyMuPDF is hosted on GitHub in this [repository](#). Because we rely on MuPDF’s C library, installation consists of two separate steps, which is why we cannot support Python’s standard `pip` process:

1. Installation of MuPDF: this involves downloading the source from their website and then compiling it on your machine.
2. Installation of PyMuPDF: this step is normal Python procedure. Usually you will have to adapt the `setup.py` to point to correct `include` and `lib` directories of your generated MuPDF.

For the Windows platform we have however combined these steps and offer binaries. This installation material is contained in a separate GitHub [repository](#) and obsoletes all other download and generation work. You only need to choose which Python version and bitness you want and then download the respective ZIP file (about 3 to 5 MB).

For installation details check out the respective chapter.

We also are registered on [PyPI](#).

There exist several demo and example programs in the main repository, ranging from simple code snippets to full-featured utilities, like text extraction, PDF joiners and bookmark maintenance.

Interesting **PDF manipulation and generation** functions have been added over time, including meta-data and bookmark maintenance, document restructuring, annotation / link handling and document creation.

1.1 Note on the Name `fitz`

The standard Python import statement for this library is `import fitz`. This has a historical reason:

The original rendering library for MuPDF was called `Libart`. “After Artifex Software acquired the MuPDF project, the development focus shifted on writing a new modern graphics library called `Fitz`. Fitz was originally intended as an R&D project to replace the aging Ghostscript graphics library, but has instead become the rendering engine powering MuPDF.” (Quoted from [Wikipedia](#)).

1.2 License

PyMuPDF is distributed under GNU GPL V3 or later.

MuPDF is distributed under a separate license, the **GNU AFFERO GPL V3**.

Both licenses apply, when you generate and use PyMuPDF and MuPDF.

Note: Version 3 of the GNU AFFERO GPL is a lot less restrictive than its earlier versions used to be. You should note that it is an “**Open Source License**”. For details consult this [website](#), especially when you want to create a commercial product with PyMuPDF.

1.3 Covered Version

This documentation covers PyMuPDF features included in build date **2017-07-30** and time **12:12:11**. These values are taken from variable `fitz.VerionDate`.

INSTALLATION

Installation generally encompasses downloading and generating PyMuPDF and MuPDF from sources.

This process consists of three steps described below under “**Option 1: Install from Sources**”.

If your operating system is Windows XP SP2 or higher (x86 or x64), you can perform a binary setup, detailed out under “**Option 2: Install from Binaries**”. This process is **a lot faster** and requires no compiler, no Visual Studio, no download of MuPDF, even no download of PyMuPDF. You only need to download the ZIP file from PyMuPDF-optional-material that fits your Python version.

2.1 Option 1: Install from Sources

2.1.1 Step 1: Download PyMuPDF

Download this repository and unzip / decompress it. This will give you a folder, let us call it PyFitz.

2.1.2 Step 2: Download and Generate MuPDF

Download `mupdf-x.xx-source.tar.gz` from <http://mupdf.com/downloads> and unzip / decompress it. Call the resulting folder `mupdf`. The latest MuPDF **development sources** are available on <https://github.com/ArtifexSoftware/mupdf> - this is **not** what you want here.

Make sure you download the (sub-) version for which PyMuPDF has stated its compatibility. The various Linux flavors usually have their own specific ways to support download of packages which we cannot cover here. Do not hesitate posting issues to our web site or sending an e-mail to the authors for getting support.

Put it inside PyFitz as a subdirectory for keeping everything in one place.

Controlling the Binary File Size:

Since version 1.9, MuPDF includes support for many dozens of additional, so-called NOTO (“no TOFU”) fonts for all sorts of alphabets from all over the world like Chinese, Japanese, Korean, Cyrillic, Indonesian, Chinese etc. If you accept MuPDF’s standard here, the resulting binary for PyMuPDF will be very big and easily approach or exceed 20 MB. The features actually needed by PyMuPDF in contrast only represent a fraction of this size: no more than 5 MB currently.

To cut off unneeded stuff from your MuPDF version, modify file `/include/mupdf/config.h` as follows:

```
#ifndef FZ_CONFIG_H

#define FZ_CONFIG_H

/ *
    Choose which plotters we need.
    By default we build the greyscale, RGB and CMYK plotters in,
    but omit the arbitrary plotters. To avoid building
    plotters in that aren't needed, define the unwanted
```

```
FZ_PLOTTERS... define to 0.
*/
/* #define FZ_PLOTTERS_G 1 */
/* #define FZ_PLOTTERS_RGB 1 */
/* #define FZ_PLOTTERS_CMYK 1 */
/* #define FZ_PLOTTERS_N 0 */

/*
    Choose which document agents to include.
    By default all but GPRF are enabled. To avoid building unwanted
    ones, define FZ_ENABLE... to 0.
*/
/* #define FZ_ENABLE_PDF 1 */
/* #define FZ_ENABLE_XPS 1 */
/* #define FZ_ENABLE_SVG 1 */
/* #define FZ_ENABLE_CBZ 1 */
/* #define FZ_ENABLE_IMG 1 */
/* #define FZ_ENABLE_TIFF 1 */
/* #define FZ_ENABLE_HTML 1 */
/* #define FZ_ENABLE_EPUB 1 */
/* #define FZ_ENABLE_GPRF 1 */

/*
    Choose whether to enable JavaScript.
    By default JavaScript is enabled both for mutool and PDF interactivity.
*/
// #define FZ_ENABLE_JS 1

/*
    Choose which fonts to include.
    By default we include the base 14 PDF fonts,
    DroidSansFallback from Android for CJK, and
    Charis SIL from SIL for epub/html.
    Enable the following defines to AVOID including
    unwanted fonts.
*/
/* To avoid all noto fonts except CJK, enable: */
#define TOFU // <===== PyMuPDF

/* To skip the CJK font, enable: */
#define TOFU_CJK // <===== PyMuPDF

/* To skip CJK Extension A, enable: */
#define TOFU_CJK_EXT // <===== PyMuPDF

/* To skip the Emoji font, enable: */
#define TOFU_EMOJI // <===== PyMuPDF

/* To skip the ancient/historic scripts, enable: */
#define TOFU_HISTORIC // <===== PyMuPDF

/* To skip the symbol font, enable: */
/* #define TOFU_SYMBOL */

/* To skip the SIL fonts, enable: */
#define TOFU_SIL // <===== PyMuPDF

/* To skip the Base14 fonts, enable: */
/* #define TOFU_BASE14 */
/* (You probably really don't want to do that except for measurement purposes!) */

/* ----- DO NOT EDIT ANYTHING UNDER THIS LINE ----- */
```

```
... ..
#endif /* FZ_CONFIG_H */
```

The above choice should bring down your binary file size to around 5 MB or less, depending on your bitness.

Generate MuPDF now.

The MuPDF source includes generation procedures / makefiles for numerous platforms. For Windows platforms, Visual Studio solution and project definitions are provided.

Consult additional installation hints on PyMuPDF's [main page](#) on Github.com. Among other things you will find a Wiki pages with details on building the Windows binaries or user provided installation experiences.

2.1.3 Step 3: Build / Setup PyMuPDF

Adjust the `setup.py` script as necessary. E.g. make sure that

- the include directory is correctly set in sync with your directory structure
- the object code libraries are correctly defined

Now perform a `python setup.py install`.

2.1.4 Using UPX

Your PyMuPDF installation will end up with four files: `__init__.py`, `fitz.py`, `utils.py` and the binary file `_fitz.xxx` in the `site-packages` directory. The extension of the binary will be `.pyd` on Windows and `.so` on other platforms.

Depending on your OS, your compiler and your font support choice (see above), this binary can be quite large and range from 5 MB to 20 MB. You can reduce this by applying the compression utility [UPX](#) to it, which exists for many operating systems. UPX will reduce the size of `_fitz.xxx` by more than 50%. You will end up with 2.5 MB to 9 MB without impacting functionality nor execution speed.

2.2 Option 2: Install from Binaries

This installation option is based on pre-built binaries for Python versions on Windows XP, 7, 8 and 10 (32bit or 64bit). Supported Python versions include 2.7 and 3.1 through 3.6 (32bit and 64bit).

2.2.1 Step 1: Download Optional Material

Download [PyMuPDF-optional-material](#). From directory `binary_setups` select the zip file corresponding to your configuration and unzip it anywhere you like. To reduce download time, only download the zip file corresponding to your Python version - a matter of less than 3 MB.

2.2.2 Step 2: Install PyMuPDF

Open a command prompt at the unzipped folder's directory that contains `setup.py` and enter `python setup.py install` (or `py setup.py install` if you have the Python launcher, see below).

You are done within 2 seconds.

This process requires no download of anything else, no compiler, no Visual Studio and is **very** fast. The only pre-requisite is, that your Python configuration matches the zip file.

2.2.3 MD5 Checksums

Binary download setup scripts contain an integrity check based on MD5 check sums.

The directory structure of each zip file `pymupdf-1.10.?.?-py??-x???.zip` is as follows:

```
fitz
├── fitz
│   ├── __init__.py
│   ├── _fitz.pyd
│   ├── fitz.py
│   └── utils.py
├── MANIFEST
├── md5.txt
├── PKG-INFO
└── setup.py
```

During setup, the MD5 check sum of the four installation files `__init__.py`, `_fitz.pyd`, `utils.py` and `fitz.py` is being calculated and compared against the pre-calculated check sum contained in file `md5.txt`. If a mismatch is detected, the error message

```
md5 mismatch:  probable download error
```

is issued and setup is cancelled. In this case, please check your download for any problems.

2.2.4 Targeting Parallel Python Installations

Setup scripts for binary install support the Python launcher `py.exe` introduced with version 3.3.

They contain **shebang lines** that specify the intended Python version, and additional checks for detecting error situations.

This can be used to target the right Python version if you have several installed in parallel (and of course the Python launcher, too). Use the following statement to set up PyMuPDF correctly:

```
py setup.py install
```

The shebang line of `setup.py` will be interpreted by `py.exe` to automatically find the right Python, and the internal checks will make sure that version and bitness are what they should be.

TUTORIAL

This tutorial will show you the use of MuPDF in Python step by step.

Because MuPDF supports not only PDF, but also XPS, OpenXPS, CBZ and EPUB formats, so does PyMuPDF. Nevertheless we will only talk about PDF files for the sake of brevity. At places where indeed only PDF files are supported, this will be mentioned explicitly.

3.1 Importing the Bindings

The Python bindings to MuPDF are made available by this import statement:

```
import fitz
```

You can check your version by printing the docstring:

```
>>> print(fitz.__doc__)
PyMuPDF 1.9.1: Python bindings for the MuPDF 1.9a library,
built on 2016-07-01 13:06:02
>>>
```

3.2 Opening a Document

In order to access a supported document, it must be opened with the following statement:

```
doc = fitz.open(filename)      # or fitz.Document(filename)
```

This will create `doc` as a *Document* object. `filename` must be a Python string or unicode object that specifies the name of an existing file.

It is also possible to open a document from memory data, i.e. without using a file, or create a new, empty PDF. See *Document* for details.

A document contains many attributes and functions. Among them are meta information (like “author” or “subject”), number of total pages, outline and encryption information.

3.3 Some Document Methods and Attributes

| Method / Attribute | Description |
|----------------------------|------------------------------|
| <i>Document.pageCount</i> | number of pages (int). |
| <i>Document.metadata</i> | metadata (dictionary). |
| <i>Document.getToC()</i> | table of contents (list). |
| <i>Document.loadPage()</i> | create a Page object. |

3.4 Accessing Meta Data

`Document.metadata` is a Python dictionary with the following keys. For details of their meanings and formats consult the PDF manuals, e.g. *Adobe PDF Reference 1.7*. Further information can also be found in chapter *Document*. The meta data fields are strings (or `None`) if not otherwise indicated. Be aware that not all of them necessarily contain meaningful data.

| Key | Value |
|--------------|-------------------------------|
| producer | producer (producing software) |
| format | PDF format, e.g. 'PDF-1.4' |
| encryption | encryption method used |
| author | author |
| modDate | date of last modification |
| keywords | keywords |
| title | title |
| creationDate | date of creation |
| creator | creating application |
| subject | subject |

3.5 Working with Outlines

The easiest way to get all outlines of a document, is creating a table of contents:

```
toc = doc.getToC()
```

This will return a Python list of lists `[[lvl, title, page, ...], ...]`.

`lvl` is the hierarchy level of the entry (starting from 1), `title` is the entry's title, and `page` the page number (1-based!). Other parameters describe details of the bookmark target.

3.6 Working with Pages

Tasks that can be performed with a *Page* are at the core of MuPDF's functionality. Among other things, you can render a page, optionally zooming, rotating, shifting or shearing it. You can write it's image to files, extract text from it or search for text strings.

At first, a page object must be created:

```
page = doc.loadPage(n)          # represents page n of the document (0-based)
page = doc[n]                   # short form
```

The integer `n` above may be any number less than the total number of pages of the document. All negative values are allowed, e.g. `doc[-1]` means the last page, as with Python lists. `doc[-500]` is **always** valid for any document: to access the respective actual page, the total number of pages is added to -500 until the result is no longer negative.

Some typical uses of *Page* objects follow:

3.6.1 Inspecting the Links of a Page

Here is how to get all links and their types:

```
# get all links of the current page
links = page.getLinks()
```

`links` is a Python list containing Python dictionaries as entries. For details see `Page.getLinks()`.

3.6.2 Rendering a Page

This example creates an image out of a page's content (default parameters shown):

```
pix = page.getPixmap(matrix = fitz.Identity,
                    colorspace = "rgb",
                    alpha = True)
```

Now `pix` contains an RGB image of the page, ready to be used. The above method offers lots of variations for increasing image precision, colorspace selection, transparency exclusion, rotation, mirroring, shifting, shearing, etc.

3.6.3 Saving the Page Image in a File

We can simply store the image in a PNG file:

```
pix.writePNG("test.png")
```

3.6.4 Displaying the Image in Dialog Managers

We can also use the image in a dialog. `Pixmap.samples` represents the area of bytes of all the pixels as a Python bytes object. This area is directly usable by presumably most dialog managers. Here are two examples. Please also have a look at the examples directory of this repository.

wxPython:

```
bitmap = wx.BitmapFromBufferRGBA(pix.width, # image width
                                pix.height,  # image height
                                pix.samples)  # bytes with pixel data
```

Tkinter:

```
# the following requires: "from PIL import Image, ImageTk"
img = Image.frombytes("RGBA", [pix.width, pix.height], pix.samples)
photo = ImageTk.PhotoImage(img)
```

Now, `photo` can be used as an image in TK.

3.6.5 Extracting Text

We can also extract all text of a page in one chunk of string:

```
text = page.getText("text")
```

For the parameter, the following values can be specified:

- `text`: plain text with line breaks (default). No format and no position info.
- `html`: line breaks, alignment, grouping in HTML syntax. No format and no position info.
- `json`: full formatting info in JSON format (except colors and fonts) down to spans (see Appendix 2). Use a `json` module to interpret.
- `xml`: full (except colors) formatting info in XML format down to each single character (!). Use an XML module to interpret.

To give you an idea about the output of these alternatives, we did text example extracts. See the Appendix 2.

3.6.6 Searching Text

You can find out, exactly where on a page a certain string appears like this:

```
>>> areas = page.searchFor("mupdf", hit_max = 16)
```

The variable `areas` will contain a list of up to 16 *Rect* rectangles, each of which surrounds one occurrence of string “mupdf” (case insensitive).

Please also do have a look at chapter *Working together: Device, DisplayList, TextPage and TextSheet* and at demo program `demo.py`. Among other things they contain details on how the *TextPage*, *TextSheet*, *Device* and *DisplayList* classes can be used for a more direct control, e.g. when performance considerations suggest it.

3.7 PDF Maintenance

Since version 1.9, PyMuPDF provides several options to modify PDF documents (only).

The `Document.save()` method automatically stores a document in its current (potentially modified) state on disk.

Be aware that a PDF document can be modified unnoticed by the user in two ways:

- During open, integrity checks are used to determine the health of the PDF structure. Any errors will automatically be corrected to present a repaired document in memory for further processing. If this is the case, the document is regarded as being modified.
- After a document has been decrypted, the document in memory obviously has changed and also counts as being modified.

In these two cases, the save method will store a repaired and / or decrypted version, and saving **must occur to a new file**.

The following describe some more intentional ways to manipulate PDF documents. Beyond those mentioned here, you can also modify the table of contents and meta information.

3.7.1 Modifying, Creating, Re-arranging and Deleting Pages

There are several ways to manipulate the page tree of a PDF:

Methods `Document.deletePage()` and `Document.deletePageRange()` delete a page (range) specified by zero-based number(s).

Methods `Document.copyPage()` and `Document.movePage()` copy or move a page to another location of the document.

`Document.insertPage()` inserts a new page, optionally containing some plain text.

Method `Document.select()` accepts a list of integers as argument. These integers must be in range $0 \leq i < \text{pageCount}$. When executed, all pages **not occurring** in this list will be deleted. Only pages that do occur will remain - **in the sequence specified and as many times as specified**.

So you can easily create new PDFs with the first or last 10 pages, only the odd or only the even pages (for doing double-sided printing), pages that **do** or **do not** contain a certain text, ... whatever you may think of.

The saved new document will contain all still valid links, annotations and bookmarks.

Pages can also be modified, deleted, copied or moved individually with a range of methods (e.g. annotation and link maintenance, text and image insertion).

3.7.2 Joining and Splitting PDF Documents

Method `Document.insertPDF()` inserts another PDF document at a specified place of the current one. Here is a simple example (`doc1` and `doc2` are opened PDF documents):

```
>>> # append complete doc2 to the end of doc1
>>> doc1.insertPDF(doc2)
```

Here is how to split `doc1`. This creates a new document of its first and last 10 pages:

```
>>> doc2 = fitz.open()
>>> doc2.insertPDF(doc1, to_page = 9)
>>> doc2.insertPDF(doc1, from_page = len(doc1) - 10)
>>> doc2.save(...)
```

More can be found in the `Document` chapter. Also have a look at `PDFjoiner.py` in the repository's `example` directory.

3.7.3 Saving

As mentioned above, `save()` will automatically **always** save the document in its current state, decrypted and / or repaired, and including all of your changes. The method's parameters offer you additional ways to (de-) compress or clean content and much more.

Since MuPDF 1.9, you can also write changes back to the original file by specifying `incremental = True`. This process is (usually) **extremely fast**, since changes are **appended to the original file** - it will not be rewritten as a whole.

`Document.save()` supports all options of MuPDF's command line utility `mutool clean`, see the following table (corresponding `mutool clean` option is indicated as "mco").

| Option | mco | Effect |
|------------------------------|-------|---|
| <code>garbage = 1</code> | -g | garbage collect unused objects |
| <code>garbage = 2</code> | -gg | in addition to 1, compact xref tables |
| <code>garbage = 3</code> | -ggg | in addition to 2, merge duplicate objects |
| <code>garbage = 4</code> | -gggg | in addition to 3, check for duplicate streams |
| <code>clean = 1</code> | -s | clean content streams |
| <code>deflate = 1</code> | -z | deflate uncompressed streams |
| <code>ascii = 1</code> | -a | convert data to ASCII format |
| <code>linear = 1</code> | -l | create a linearized version (do not use yet) |
| <code>expand = 1</code> | -i | decompress images |
| <code>expand = 2</code> | -f | decompress fonts |
| <code>expand = 255</code> | -d | decompress all |
| <code>incremental = 1</code> | n/a | append changes to the original |

Be ready to experiment a little if you want to fully exploit above options: like with `mutool clean`, not all combinations may always work: there are just too many ill-constructed PDF files out there ...

We have found, that the combination `mutool clean -gggg -z` yields excellent compression results and is very stable. In PyMuPDF this corresponds to `doc.save(filename, garbage=4, deflate=1)`.

3.8 Closing

It is often desirable to "close" a document to relinquish control of the underlying file to the OS, while your program is still running.

This can be achieved by the `Document.close()` method. Apart from closing the underlying file, buffer areas associated with the document will be freed (if the document has been created from memory data, only the buffer release will take place).

3.9 Example: Dynamically Cleaning up Corrupt PDF Documents

This shows a potential use of PyMuPDF with another Python PDF library (`pdfrw`).

If a clean, non-corrupt or decompressed PDF is needed, one could dynamically invoke PyMuPDF to recover from problems like so:

```
import sys
from pdfrw import PdfReader
import fitz
from io import BytesIO

#-----
# 'tolerant' PDF reader
#-----
def reader(fname):
    ifile = open(fname, "rb")
    idata = ifile.read()           # put in memory
    ifile.close()
    ibuffer = BytesIO(idata)      # convert to stream
    try:
        return PdfReader(ibuffer) # let us try
    except:
        # problem! heal it with PyMuPDF
        doc = fitz.open("pdf", idata) # open and save a corrected
        c = doc.write(garbage = 4)    # version in memory
        doc.close()
        doc = idata = None            # free storage
        ibuffer = BytesIO(c)          # convert to stream
        return PdfReader(ibuffer)     # let pdfrw retry
#-----
# Main program
#-----
pdf = reader("pymupdf.pdf")
print pdf.Info
# do further processing
```

With the command line utility `pdftk` (available for Windows only) a similar result can be achieved, see [here](#). However, you must invoke it as a separate process via `subprocess.Popen`, using `stdin` and `stdout` as communication vehicles.

CLASSES

4.1 Annot

Quote from the *Adobe PDF Reference 1.7*: “An annotation associates an object such as a note, sound, or movie with a location on a page of a PDF document, or provides a way to interact with the user by means of the mouse and keyboard.”

This class supports accessing such annotations - not only for PDF files, but for all MuPDF supported document types. However, only a few methods and properties apply to non-PDF documents.

There is a parent-child relationship between an annotation and its page. If the page object becomes unusable (closed document, any document structure change, etc.), then so does every of its existing annotation objects - an exception is raised saying that the object is “orphaned”, whenever an annotation property or method is accessed.

| Attribute | Short Description |
|----------------------------|---|
| <i>Annot.getPixmap()</i> | image of the annotation as a pixmap |
| <i>Annot.setInfo()</i> | PDF only: change metadata of an annotation |
| <i>Annot.setBorder()</i> | PDF only: changes the border of an annotation |
| <i>Annot.setFlags()</i> | PDF only: changes the flags of an annotation |
| <i>Annot.setRect()</i> | PDF only: changes the rectangle of an annotation |
| <i>Annot.setColors()</i> | PDF only: changes the colors of an annotation |
| <i>Annot.updateImage()</i> | PDF only: applies border and color values to shown image |
| <i>Annot.fileInfo()</i> | PDF only: returns attached file information |
| <i>Annot.fileGet()</i> | PDF only: returns attached file content |
| <i>Annot.fileUpd()</i> | PDF only: sets attached file new content |
| <i>Annot.border</i> | PDF only: border details |
| <i>Annot.colors</i> | PDF only: border / background and fill colors |
| <i>Annot.flags</i> | PDF only: annotation flags |
| <i>Annot.info</i> | PDF only: various information |
| <i>Annot.lineEnds</i> | PDF only: start / end appearance of line-type annotations |
| <i>Annot.next</i> | link to the next annotation |
| <i>Annot.parent</i> | page object of the annotation |
| <i>Annot.rect</i> | rectangle containing the annotation |
| <i>Annot.type</i> | PDF only: type of the annotation |
| <i>Annot.vertices</i> | PDF only: point coordinates of Polygons, PolyLines, etc. |

Class API

class Annot

`getPixmap(matrix = fitz.Identity, colorspace = fitz.csRGB, alpha = False)`

Creates a pixmap from the annotation as it appears on the page in untransformed coordinates.
The pixmap's *IRect* equals *Annot.rect.irect* (see below).

Parameters

- **matrix** (*Matrix*) – a matrix to be used for image creation. Default is the `fitz.Identity` matrix.
- **colorspace** (*Colorspace*) – a colorspace to be used for image creation. Default is `fitz.csRGB`.
- **alpha** (*bool*) – whether to include transparency information. Default is `False`.

Return type *Pixmap*

setInfo(*d*)

Changes the info dictionary. This includes dates, contents, subject and author (title). Changes for **name** will be ignored.

Parameters *d* (*dict*) – a dictionary compatible with the **info** property (see below). All entries must be `unicode`, `bytes`, or strings. If `bytes` values in Python 3 they will be treated as being UTF8 encoded.

setRect(*rect*)

Changes the rectangle of an annotation. The annotation can be moved around and both sides of the rectangle can be independently scaled. However, the annotation appearance will never get rotated, flipped or sheared.

Parameters *rect* (*Rect*) – the new rectangle of the annotation. This could e.g. be a rectangle `rect = Annot.rect * M` with a suitable *Matrix* *M* (only scaling and translating will yield the expected effect).

setBorder(*value*)

PDF only: Change border width and dashing properties. Any other border properties will be deleted.

Parameters *value* (*float or dict*) – a number or a dictionary specifying the desired border properties. If a dictionary, its **width** and **dashes** keys are used (see property **annot.border**). If a number is specified or a dictionary like `{"width": w}`, only the border width will be changed and any dashes will remain unchanged. Conversely, with a dictionary `{"dashes": [...]}`, only line dashing will be changed. To remove dashing and get a contiguous line, specify `{"dashes": []}`.

setFlags(*flags*)

Changes the flags of the annotation. See *Annotation Flags* for possible values and use the `|` operator to combine several.

Parameters *flags* (*int*) – an integer specifying the required flags.

setColors(*d*)

Changes the colors associated with the annotation.

Parameters *d* (*dict*) – a dictionary containing color specifications. For accepted dictionary keys and values see below. The most practical way should be to first make a copy of the **colors** property and then modify this dictionary as required.

updateImage()

Attempts to modify the displayed graphical image such that it coincides with the values currently contained in the **border** and **colors** properties. This is achieved by modifying the contents stream of the associated appearance *XObject*. Not all possible formats of content streams are currently supported: if the stream contains invocations of yet other *XObject* objects, a `ValueError` is raised.

fileInfo()

Returns basic information of an attached file (file attachment annotations only).

Return type *dict*

Returns a dictionary with keys **filename**, **size** (uncompressed file size), **length** (compressed length).

`fileGet()`

Returns the uncompressed content of the attached file.

Return type bytes or str (Py2)

Returns the content of the attached file.

`fileUpd(buffer, filename=None)`

Updates the content of an attached file with new data. Optionally, the filename can be changed, too.

Parameters

- `buffer` (*bytes or bytearray*) – the new file content.
- `filename` (*str*) – new filename to associate with the file.

Return type int

Returns zero

`parent`

The owning page object of the annotation.

Return type *Page*

`rect`

The rectangle containing the annotation in untransformed coordinates.

Return type *Rect*

`next`

The next annotation on this page or None.

Return type *Annot*

`type`

Meaningful for PDF only: A number and one or two strings describing the annotation type, like [2, 'FreeText', 'FreeTextCallout']. The second string entry is optional and may be empty. [] if not PDF. See the appendix *Annotation Types* for a list of possible values and their meanings.

Return type list

`info`

Meaningful for PDF only: A dictionary containing various information. All fields are unicode or strings (Python 2 or Python 3 respectively).

- `name` - e.g. for [12, 'Stamp'] type annotations it will contain the stamp text like `Sold` or `Experimental`.
- `content` - a string containing the text for type `Text` and `FreeText` annotations. Commonly used for filling the text field of annotation pop-up windows. For `FileAttachment` it should be used as description for the attached file. Initially just contains the filename.
- `title` - a string containing the title of the annotation pop-up window. By convention, this is used for the annotation author.
- `creationDate` - creation timestamp.
- `modDate` - last modified timestamp.
- `subject` - subject, an optional string.

Return type dict

`flags`

Meaningful for PDF only: An integer whose low order bits contain flags for how the annotation should be presented. See section *Annotation Flags* for details.

Return type int

lineEnds

Meaningful for PDF only: A dictionary specifying the starting and the ending appearance of annotations of types `Line`, `PolyLine`, among others. An example would be `{'start': 'None', 'end': 'OpenArrow'}`. `{}` if not specified or not applicable. For possible values and descriptions in this list, see the [Adobe PDF Reference 1.7](#), table 8.27 on page 630.

Return type dict

vertices

Meaningful for PDF only: A list containing point (“vertices”) coordinates (each given by 2 floats specifying the x and y coordinate respectively) for various types of annotations:

- `Line` - the starting and ending coordinates (4 floats).
- `[2, 'FreeText', 'FreeTextCallout']` - 4 or 6 floats designating the starting, the (optional) knee point, and the ending coordinates.
- `PolyLine` / `Polygon` - the coordinates of the edges connected by line pieces ($2 * n$ floats for n points).
- text markup annotations - $8 * n$ floats specifying the `QuadPoints` of the n marked text spans (see [Adobe PDF Reference 1.7](#), page 634).
- `Ink` - list of one to many sublists of vertex coordinates. Each such sublist represents a separate line in the drawing.

Return type list

colors

Meaningful for PDF only: A dictionary of two lists of floats in range $0 \leq \text{float} \leq 1$ specifying the common (`common`) or `stroke` and the interior (`fill`) **non-stroke** colors. The common color is used for borders and everything that is actively painted or written (“*stroked*”). The fill color is used for the interior of objects like line ends, circles and squares. The lengths of these lists implicitly determine the colorspaces used: 1 = GRAY, 3 = RGB, 4 = CMYK. So `[1.0, 0.0, 0.0]` stands for RGB and color `red`. Both lists can be `[]` if not specified. The dictionary will be empty `{}` if no PDF. The value of each float is mapped to integer values from 0 ($\Leftrightarrow 0.0$) to 255 ($\Leftrightarrow 1.0$).

Return type dict

border

Meaningful for PDF only: A dictionary containing border characteristics. It will be empty `{}` if not PDF or when no border information is provided. Technically, the PDF entries `/Border`, `/BS` and `/BE` will be checked to build this information. The following keys can occur:

- `width` - a float indicating the border thickness in points.
- `effect` - a list specifying a border line effect like `[1, 'C']`. The first entry “intensity” is an integer (from 0 to 2 for maximum intensity). The second is either ‘S’ for “no effect” or ‘C’ for a “cloudy” line.
- `dashes` - a list of integers (arbitrarily limited to 10) specifying a line dash pattern in user units (usually points). `[]` means no dashes, `[n]` means equal on-off lengths of n points, longer lists will be interpreted as specifying alternating on-off length values. See the [Adobe PDF Reference 1.7](#) page 217 for more details.
- `style` - 1-byte border style: S (Solid) = solid rectangle surrounding the annotation, D (Dashed) = dashed rectangle surrounding the annotation, the dash pattern is specified by the `dashes` entry, B (Beveled) = a simulated embossed rectangle that appears to be raised above the surface of the page, I (Inset) = a simulated engraved rectangle that appears to be recessed below the surface of the page, U (Underline) = a single line along the bottom of the annotation rectangle.

Return type dict

4.1.1 Example

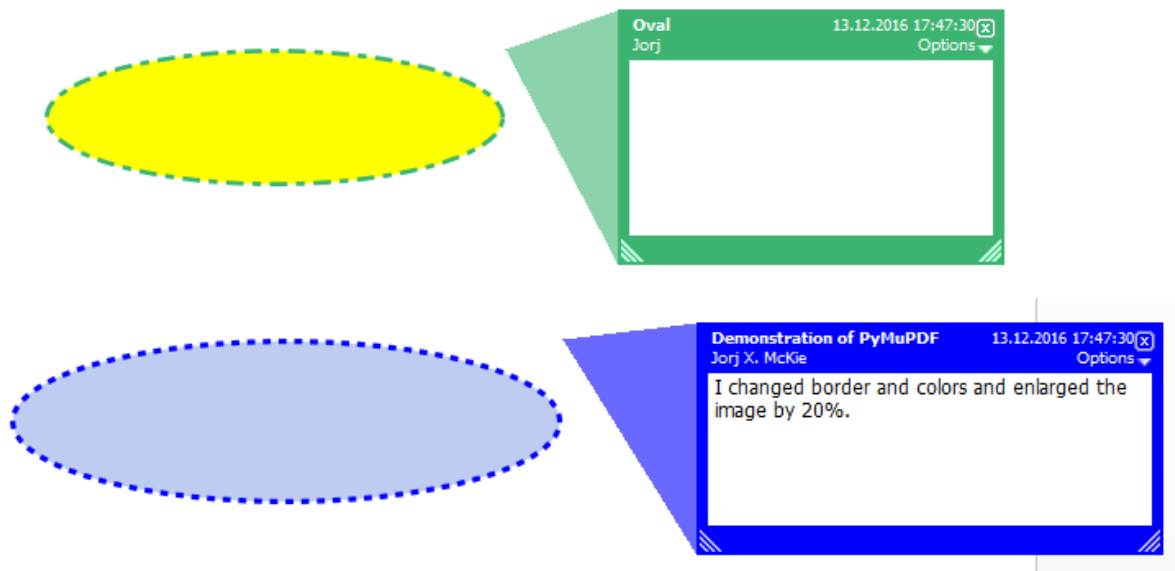
Change the graphical image of an annotation. Also update the “author” and the text to be shown in the popup window:

```
doc = fitz.open("circle-in.pdf")
page = doc[0]
annot = page.firstAnnot
annot.setBorder({"dashes": [3]})

# set border / popup color to blue and fill color to some light blue
annot.setColors({"common": [0, 0, 1], "fill": [0.75, 0.8, 0.95]})
info = annot.info
info["title"] = "Jorj X. McKie"

# text in popup window ...
info["content"] = "I changed border and colors and enlarged the image by 20%."
info["subject"] = "Demonstration of PyMuPDF"
annot.setInfo(info)
r = annot.rect
r.x1 = r.x0 + r.width * 1.2
r.y1 = r.y0 + r.height * 1.2
annot.setRect(r)
annot.updateImage()
doc.save("circle-out.pdf", garbage=4)
```

This is how the circle annotation looks like, before and after the change:



4.2 Colorspace

Represents the color space of a *Pixmap*.

Class API

```
class Colorspace
```

`--init__(self, n)`
Constructor

Parameters `n` (*int*) – A number identifying the colorspace. Possible values are `CS_RGB`, `CS_GRAY` and `CS_CMYK`.

name
The name identifying the colorspace. Example: `fitz.csCMYK.name = 'DeviceCMYK'`.

Type `str`

n

The number of bytes required to define the color of one pixel. Example: `fitz.csCMYK.n = 4`.

type `int`

Predefined Colorspaces

For saving some typing effort, there exist predefined colorspace objects for the three available cases.

- `csRGB = fitz.Colorspace(fitz.CS_RGB)`
- `csGRAY = fitz.Colorspace(fitz.CS_GRAY)`
- `csCMYK = fitz.Colorspace(fitz.CS_CMYK)`

4.3 Document

This class represents a document. It can be constructed from a file or from memory.

Since version 1.9.0 there exists the alias `open` for this class.

For additional details on **embedded files** refer to Appendix 3.

| Method / Attribute | Short Description |
|---|---|
| <code>Document.authenticate()</code> | decrypt the document |
| <code>Document.close()</code> | close the document |
| <code>Document.copyPage()</code> | PDF only: copy a page to another location |
| <code>Document.deletePage()</code> | PDF only: delete a page by its number |
| <code>Document.deletePageRange()</code> | PDF only: delete a range of pages |
| <code>Document.embeddedFileAdd()</code> | PDF only: add a new embedded file from buffer |
| <code>Document.embeddedFileDel()</code> | PDF only: delete an embedded file entry |
| <code>Document.embeddedFileGet()</code> | PDF only: extract an embedded file buffer |
| <code>Document.embeddedFileInfo()</code> | PDF only: metadata of an embedded file |
| <code>Document.embeddedFileSetInfo()</code> | PDF only: change metadata of an embedded file |
| <code>Document.getPageFontList()</code> | make a list of fonts on a page |
| <code>Document.getPageImageList()</code> | make a list of images on a page |
| <code>Document.getPagePixmap()</code> | create a pixmap of a page by page number |
| <code>Document.getPageText()</code> | extract the text of a page by page number |
| <code>Document.getToC()</code> | create a table of contents |
| <code>Document.insertPage()</code> | PDF only: insert a new empty page |
| <code>Document.insertPDF()</code> | PDF only: insert pages from another PDF |
| <code>Document.loadPage()</code> | read a page |
| <code>Document.movePage()</code> | PDF only: move a page to another location |
| <code>Document.save()</code> | PDF only: save the document |
| <code>Document.saveIncr()</code> | PDF only: save the document incrementally |
| <code>Document.searchPageFor()</code> | search for a string on a page |
| <code>Document.select()</code> | PDF only: select a subset of pages |

Continued on next page

Table 4.1 – continued from previous page

| Method / Attribute | Short Description |
|---|---|
| <code>Document.setMetadata()</code> | PDF only: set the metadata |
| <code>Document.setToC()</code> | PDF only: set the table of contents (TOC) |
| <code>Document.write()</code> | PDF only: writes the document to memory |
| <code>Document.embeddedFileCount</code> | number of embedded files |
| <code>Document.isClosed</code> | has document been closed? |
| <code>Document.metadata</code> | metadata |
| <code>Document.name</code> | filename of document |
| <code>Document.needsPass</code> | require password to access data? |
| <code>Document.openErrCode</code> | > 0 if repair occurred during open |
| <code>Document.openErrMsg</code> | last error message if openErrCode > 0 |
| <code>Document.outline</code> | first <i>Outline</i> item |
| <code>Document.pageCount</code> | number of pages |
| <code>Document.permissions</code> | permissions to access the document |

Class API

class Document

`--init__(self[, filename])`

Constructs a Document object from filename.

Parameters `filename` (*str*) – A string containing the path / name of the document file to be used. The file will be opened and remain open until either explicitly closed (see below) or until end of program. If omitted or `None`, a new empty PDF document will be created.

Return type Document

Returns A Document object.

`--init__(self, filetype, stream)`

Constructs a Document object from memory stream.

Parameters

- `filetype` (*str*) – A string specifying the type of document contained in `stream`. This may be either something that looks like a filename (e.g. "x.pdf"), in which case MuPDF uses the extension to determine the type, or a mime type like `application/pdf`. Recommended is using the filename scheme, or even the name of the original file for documentation purposes. But just using strings like "pdf" will also work.
- `stream` (*bytearray or bytes or (Python 2 only) str*) – A memory area representing the content of a supported document type.

Return type Document

Returns A Document object.

`authenticate(password)`

Decrypts the document with the string `password`. If successful, all of the document's data can be accessed (e.g. for rendering).

Parameters `password` (*str*) – The password to be used.

Return type int

Returns `True` (1) if decryption with `password` was successful, `False` (0) otherwise. If successful, indicator `isEncrypted` is set to `False`.

`loadPage(pno)`

Loads a *Page* for further processing like rendering, text searching, etc. See the *Page* object.

Parameters `pno` (*int*) – page number, zero-based (0 is the first page of the document) and `< doc.pageCount`. If `number < 0`, then `page number % pageCount` will be loaded (IAW `pageCount` will be added to `number` repeatedly, until the result is no longer negative). For example: in order to load the last page, you can specify `doc.loadPage(-1)`. After this you have `page.number == doc.pageCount - 1`.

Return type *Page*

Note: Conveniently, pages can also be loaded via indexes over the document: `doc.loadPage(n) == doc[n]`. Consequently, a document can also be used as an iterator over its pages, e.g. `for page in doc: ...` and `for page in reversed(doc): ...` will yield the *Page* objects of `doc` as `page`.

`getToC(simple = True)`

Creates a table of contents out of the document's outline chain.

Parameters `simple` (*bool*) – Indicates whether a detailed ToC is required. If `simple == False`, each entry of the list also contains a dictionary with *linkDest* details for each outline entry.

Return type list

Returns a list of lists. Each entry has the form `[lvl, title, page, dest]`. Its entries have the following meanings:

- `lvl` - hierarchy level (integer). The first entry has hierarchy level 1, and entries in a row increase by at most one level.
- `title` - title (string)
- `page` - 1-based page number (integer). Page numbers `< 1` either indicate a target outside this document or no target at all (see next entry).
- `dest` - included only if `simple = False` is specified. A dictionary containing details of the link destination.

`getPagePixmap(pno, matrix = fitz.Identity, colorspace = "rgb", clip = None, alpha = True)`

Creates a pixmap from page `pno` (zero-based).

Parameters

- `pno` (*int*) – Page number, zero-based. Any value `< len(doc)` is acceptable.
- `matrix` (*Matrix*) – A transformation matrix - default is *Identity*.
- `colorspace` (str or *Colorspace*) – A string specifying the requested colorspace - default is `rgb`.
- `clip` (*IRect*) – An *IRect* to restrict rendering of the page to the rectangle's area. If not specified, the complete page will be rendered.
- `alpha` (*bool*) – Indicates whether transparency should be included. Choose `False` if not absolutely required, as it saves memory considerably (25% for RGB).

Return type *Pixmap*

`getPageImageList(pno)`

PDF only: Returns a nested list of all image descriptions referenced by a page.

Parameters `pno` (*int*) – page number, zero-based. Any value `< len(doc)` is acceptable.

Return type list

Returns

a list of images shown on this page. Each entry looks like `[xref, gen, width, height, bpc, colorspace, alt. colorspace, name]`. Where `xref` is the image object number, `gen` its generation number (should usually be zero), `width` and `height` are the image dimensions, `bpc` denotes the number of bits per component (a typical value is 8), `colorspace` a string naming the colorspace (like `DeviceRGB`), `alt. colorspace` is any alternate colorspace depending on the value of `colorspace`, and `name` - which is the symbolic name (*str*) by which the page references this particular image in its content stream. See below how this information can be used to extract pages images as separate files. Another demonstration:

```
>>> doc = fitz.open("pymupdf.pdf")
>>> imglist = doc.getPageImageList(0)
>>> for img in imglist: print img
[[241, 0, 1043, 457, 8, 'DeviceRGB', '', 'Im1']]
>>> pix = fitz.Pixmap(doc, 241)
>>> pix
fitz.Pixmap(DeviceRGB, fitz.IRect(0, 0, 1043, 457), 0)
```

getPageFontList(*pno*)

PDF only: Return a nested list of all fonts referenced by the page.

Parameters *pno* (*int*) – page number, zero-based. Any value `< len(doc)` is acceptable.

Return type list

Returns

a list of fonts referenced by this page. Each entry looks like `[xref, gen, type, basefont, name]`. Where `xref` is the image object number, `gen` its generation number (should usually be zero), `type` is the font type (like `Type1`, `TrueType`), `basefont` is the base font name und `name` is the name by which the page references the font:

```
>>> doc = fitz.open("pymupdf.pdf")
>>> fontlist = doc.getPageFontList(85)
>>> for font in fontlist: print font
[1024, 0, 'Type1', 'CJXQIC+NimbusMonL-Bold', 'R366']
[141, 0, 'Type1', 'HJQJNS+NimbusMonL-Regu', 'R247']
[162, 0, 'Type1', 'CTCORW+NimbusRomNo9L-Regu', 'R245']
[1039, 0, 'Type1', 'PWJUZZ+NimbusRomNo9L-ReguItal', 'R373']
[202, 0, 'Type1', 'VMQYGP+NimbusRomNo9L-Medi', 'R243']
[100, 0, 'Type1', 'LSBBMD+NimbusSanL-Bold', 'R201']
```

Note: A document often contain more fonts, than any single page uses. The name reference is specific for the page, i.e. other pages using the same font may use a different name for it. A font appearing in this list does not necessarily mean that there actually exists text for it on that page. But conversely, every piece of text on the page will use one of these fonts.

Note: For more background see *Adobe PDF Reference 1.7* chapters 5.4 to 5.8, pp 410.

getPageText(*pno*, *output* = "text")

Extracts the text of a page given its page number *pno* (zero-based).

Parameters

- *pno* (*int*) – Page number, zero-based. Any value `< len(doc)` is acceptable.

- **output** (*str*) – A string specifying the requested output format: text, html, json or xml. Default is **text**.

Return type str

select(*list*)

PDF only: Keeps only those pages of the document whose numbers occur in the list. Empty lists or elements outside the range $0 \leq \text{page} < \text{doc.pageCount}$ will cause a **ValueError**. For more details see remarks at the bottom of this chapter.

Parameters *list* (*sequence*) – A list (or tuple) of page numbers (zero-based) to be included. Pages not in the list will be deleted (from memory) and become unavailable until the document is reopened. **Page numbers can occur multiple times and in any order:** the resulting sub-document will reflect the list exactly as specified.

Return type int

Returns Zero upon successful execution. All document information will be updated to reflect the new state of the document, like outlines, number and sequence of pages, etc. Changes become permanent only after saving the document. Incremental save is supported.

setMetadata(*m*)

PDF only: Sets or updates the metadata of the document as specified in *m*, a Python dictionary. As with method **select()**, these changes become permanent only when you save the document. Incremental save is supported.

Parameters *m* (*dict*) – A dictionary with the same keys as **metadata** (see below). All keys are optional. A PDF's format and encryption method cannot be set or changed, these keys therefore have no effect and will be ignored. If any value should not contain data, do not specify its key or set the value to **None**. If you use *m* = {} all metadata information will be cleared to **none**. If you want to selectively change only some values, modify **doc.metadata** directly and use it as the argument for this method.

Return type int

Returns Zero upon successful execution and **doc.metadata** will be updated.

setToC(*toc*)

PDF only: Replaces the **complete current outline** tree (table of contents) with a new one. After successful execution, the new outline tree can be accessed as usual via method **getToC()** or via property **outline**. Like with other output-oriented methods, changes become permanent only via **save()** (incremental save supported). Internally, this method consists of the following two steps. For a demonstration see example below.

- Step 1 deletes all existing bookmarks.
- Step 2 creates a new TOC from the entries contained in *toc*.

Parameters *toc* (*list*) – A Python list with **all bookmark entries** that should form the new table of contents. Each entry of this list is again a list with the following format. Output variants of method **getToC()** are acceptable as input, too.

- [*lvl*, *title*, *page*, *dest*], where
- *lvl* is the hierarchy level (int > 0) of the item, starting with 1 and being at most 1 higher than that of the predecessor,
- *title* (str) is the title to be displayed.
- *page* (int) is the target page number (**attention: 1-based to support getToC()-output**), must be in valid page range if positive. Set this to -1 if there is no target, or the target is external.

- **dest** (optional) is a dictionary or a number. If a number, it will be interpreted as the desired height (in points) this entry should point to on **page** in the current document. Use a dictionary (like the one given as output by `getToC(simple = False)`) if you want to store destinations that are either “named”, or reside outside this document (other files, internet resources, etc.).

Return type int

Returns `outline` and `getToC()` will be updated upon successful execution. The return code will either equal the number of inserted items (`len(toc)`) or the number of deleted items if `toc = []`.

`save(outfile, garbage=0, clean=0, deflate=0, incremental=0, ascii=0, expand=0, linear=0)`

PDF only: Saves the document in its **current state** under the name `outfile`. A document may have changed for a number of reasons: e.g. after a successful **authenticate**, a decrypted copy will be saved, and, in addition (even without optional parameters), some basic cleaning may also have occurred, e.g. broken xref tables may have been repaired and earlier incremental changes may have been resolved. If you executed any modifying methods like `select()`, `setMetadata()`, `setToC()`, etc., their results will also be reflected in the saved version.

Parameters

- **outfile** (*str*) – The file name to save to. Must be different from the original value if `incremental=False`. When saving incrementally, **garbage** and **linear must be False / 0** and **outfile must equal** the original filename (for convenience use `doc.name`).
- **garbage** (*int*) – Do garbage collection: 0 = none, 1 = remove unused objects, 2 = in addition to 1, compact xref table, 3 = in addition to 2, merge duplicate objects, 4 = in addition to 3, check streams for duplication. Excludes **incremental**.
- **clean** (*int*) – Clean content streams¹: 0 / False, 1 / True.
- **deflate** (*int*) – Deflate uncompressed streams: 0 / False, 1 / True.
- **incremental** (*int*) – Only save changed objects: 0 / False, 1 / True. Excludes **garbage** and **linear**. Cannot be used for decrypted files and for files opened in repair mode (`openErrCode > 0`). In these cases saving to a new file is required.
- **ascii** (*int*) – Where possible make the output ASCII: 0 / False, 1 / True.
- **expand** (*int*) – Decompress contents: 0 = none, 1 = images, 2 = fonts, 255 = all. This convenience option generates a decompressed file version that can be better read by some other programs.
- **linear** (*int*) – Save a linearised version of the document: 0 = False, 1 = True. This option creates a file format for improved performance when read via internet connections. Excludes **incremental**.

Return type int

Returns Zero upon successful execution.

`saveIncr()`

PDF only: saves the document incrementally. This is a convenience abbreviation for `doc.save(doc.name, incremental = True)`.

Caution: A PDF may not be encrypted, but still be password protected against changes - see the `permissions` property. Performing incremental saves if `permissions["edit"] == False`

¹ Content streams describe what (e.g. text or images) appears where and how on a page. PDF uses a specialized language to do this (pp. 985 in *Adobe PDF Reference 1.7*), which gets interpreted when a page is loaded.

can lead to unpredictable results. Save to a new file in such a case. We also consider raising an exception under this condition.

`searchPageFor(pno, text, hit_max = 16)`

Search for `text` on page number `pno`. Works exactly like the corresponding `Page.searchFor()`. Any integer `pno < len(doc)` is acceptable.

`write(garbage=0, clean=0, deflate=0, ascii=0, expand=0, linear=0)`

PDF only: Writes the **current content of the document** to a bytearray instead of to a file like `save()`. Obviously, you should be wary about memory requirements. The meanings of the parameters exactly equal those in `Document.save()`. The tutorial contains an example for using this method as a pre-processor to `pdfw`.

Return type bytearray

Returns a bytearray containing the complete document data.

`insertPDF(docsrc, from_page = -1, to_page = -1, start_at = -1, rotate = -1, links = True)`

PDF only: Copy the page range **[from_page, to_page]** (including both) of PDF document `docsrc` into the current one. Inserts will start with page number `start_at`. Negative values can be used to indicate default values. All pages thus copied will be rotated as specified. Links can be excluded in the target, see below. All page numbers are zero-based.

Parameters

- `docsrc` (`Document`) – An opened PDF **Document** which must not be the current document object. However, it may refer to the same underlying file.
- `from_page` (`int`) – First page number in `docsrc`. Default is zero.
- `to_page` (`int`) – Last page number in `docsrc` to copy. Default is the last page.
- `start_at` (`int`) – First copied page will become page number `start_at` in the destination. If omitted, the page range will be appended to current document. If zero, the page range will be inserted before current first page.
- `rotate` (`int`) – All copied pages will be rotated by the provided value (degrees). If you do not specify a value (or -1), the original will not be changed. Otherwise it must be an integer multiple of 90 (not checked). Rotation is counter-clockwise if `rotate` is positive, else clockwise.
- `links` (`bool`) – Choose whether (internal and external) links should be included with the copy. Default is `True`. An **internal** link is always excluded if its destination is not one of the copied pages.

Return type int

Returns Zero upon successful execution.

Note: If `from_page > to_page`, pages will be copied in reverse order. If `0 <= from_page == to_page`, then one page will be copied.

Note: `docsrc` bookmarks **will not be copied**. It is easy however, to recover a table of contents for the resulting document. Look at the examples below and at program `PDFjoiner.py` in the *examples* directory: it can join PDF documents and at the same time piece together respective parts of the tables of contents.

`insertPage(to = -1, text = None, fontsize = 11, width = 595, height = 842, fontname = "Helvetica", fontfile = None, color = (0, 0, 0))`

PDF only: Insert an empty page. Default page dimensions are those of A4 portrait paper format. Optionally, text can also be inserted - provided as a string or a sequence.

Parameters

- **to** (*int*) – page number (0-based) in front of which to insert. The default -1 indicates end of document.
- **text** (*str or sequence*) – optional text to put on the page. If given, it will start at 72 points (one inch) below top and 50 points from left. Line breaks (`\n`) will be honored, if it is a string. No care will be taken as to whether lines are too wide. However, text output stops when no more lines will fit on the page (discarding any remaining text). If a sequence is specified, its entries must be a of type string. Each entry will be put on one line. Line breaks *within an entry* will be treated as any other white space. If you want to calculate the number of lines fitting on a page beforehand, use this formula: `int((height - 108) / (fontsize * 1.2))`. So, this methods reserves one inch at the top and 1/2 inches at the bottom of the page as free space.
- **fontsize** (*float*) – font size in pixels. Default is 11. If more than one line is provided, a line spacing of `fontsize * 1.2` (fontsize plus 20%) is used.
- **width** (*float*) – width in pixels. Default is 595 (A4 width). Choose 612 for *Letter width*.
- **height** (*float*) – page height in pixels. Default is 842 (A4 height). Choose 792 for *Letter height*.
- **fontname** (*str*) – name of one of the *PDF Base 14 Fonts* (default is “Helvetica”) if fontfile is not specified.
- **fontfile** (*str*) – file path of a font existing on the system. If this parameter is specified, specifying **fontname** is **mandatory**. If the font is new to the PDF, it will be embedded. Of the font file, index 0 is used. Be sure to choose a font that supports horizontal, left-to-right spacing.
- **color** (*sequence*) – RGB text color specified as a triple of floats in range 0 to 1. E.g. specify black (default) as (0, 0, 0), red as (1, 0, 0), some gray value as (0.5, 0.5, 0.5), etc.

Return type `int`

Returns number of text lines put on the page. Use this to check which part of your text did not fit.

Notes:

This method can be used to

1. create a PDF containing only one empty page of a given dimension. The size of such a file is well below 500 bytes and hence close to the theoretical PDF minimum.
2. create a protocol page of which files have been embedded, or separator pages between joined pieces of PDF Documents.
3. convert textfiles to PDF like in the demo script `text2pdf.py`.
4. For now, the inserted text should restrict itself to one byte character codes.
5. An easy way to create pages with a usual paper format, use a statement like `width, height = fitz.PaperSize("A4-L")`.
6. In order to simplify color specification, we provide a *Color Database*. This allows you to specify `color = getColor("turquoise")`, without bothering about any more details.

`deletePage(pno)`

PDF only: Delete a page given by its 0-based number in range `0 <= pno < len(doc)`.

Parameters `pno` (*int*) – the page to be deleted.

`deletePageRange(from_page = -1, to_page = -1)`

PDF only: Delete a range of pages specified as 0-based numbers. Any negative parameter will first be replaced by `len(doc) - 1`. After that, condition `0 <= from_page <= to_page < len(doc)` must be true. If the parameters are equal, one page will be deleted.

Parameters

- `from_page` (*int*) – the first page to be deleted.
- `to_page` (*int*) – the last page to be deleted.

`copyPage(pno, to = -1)`

PDF only: Copy a page within the document.

Parameters

- `pno` (*int*) – the page to be copied. Number must be in range `0 <= pno < len(doc)`.
- `to` (*int*) – the page number in front of which to insert the copy. To insert at end of document (default), specify a negative value.

`movePage(pno, to = -1)`

PDF only: Move (copy and then delete original) page to another location.

Parameters

- `pno` (*int*) – the page to be moved. Number must be in range `0 <= pno < len(doc)`.
- `to` (*int*) – the page number in front of which to insert the moved page. To insert at end of document (default), specify a negative value. Must not be in `(pno, pno + 1)`.

`embeddedFileInfo(n)`

PDF only: Retrieve information of an embedded file identified by either its number or by its name.

Parameters `n` (*int or str*) – index or name of entry. Obviously `0 <= n < embeddedFileCount` must be `True` if `n` is an integer.

Return type dict

Returns

a dictionary with the following keys:

- `name` - (*str*) name under which this entry is stored
- `file` - (*str*) filename associated with the entry
- `desc` - (*str*) description of the entry
- `size` - (*int*) original content size
- `length` - (*int*) compressed content length

`embeddedFileSetInfo(n, filename = filename, desc = desc)`

PDF only: Change some information of an embedded file given its entry number or name. At least one of `filename` and `desc` must be specified. Response will be zero if successful, else an exception is raised.

Parameters

- `n` (*int or str*) – index or name of entry. Obviously `0 <= n < embeddedFileCount` must be `True` if `n` is an integer.
- `filename` (*str*) – sets the filename of the entry.
- `desc` (*str*) – sets the description of the entry.

`embeddedFileGet(n)`

PDF only: Retrieve the content of embedded file by its entry number or name. If the document is not a PDF, or entry cannot be found, an exception is raised.

Parameters *n* (*int or str*) – index or name of entry. Obviously $0 \leq n < \text{embeddedFileCount}$ must be `True` if *n* is an integer.

Return type `bytes` (Python 3), `str` (Python 2)

`embeddedFileDel(name)`

PDF only: Remove an entry from the portfolio. As always, physical deletion of the embedded file content (and file space regain) will occur when the document is saved to a new file with `garbage` option. With an incremental save, the associated object will only be marked deleted.

Note: We do not support entry **numbers** for this function yet. If you need to e.g. delete **all** embedded files, scan through all embedded files by number, and use the returned dictionary's **name** entry to delete each one. This function will delete the first entry with this name it finds. Be wary that for arbitrary PDF files, this may not have been the only one, because PDF itself has no mechanism to prevent duplicate entries ...

Parameters *name* (*str*) – name of entry.

`embeddedFileAdd(stream, name, filename = filename, desc = desc)`

PDF only: Add new content to the document's portfolio.

Parameters

- *stream* (*bytes or bytearray or str (Python 2 only)*) – contents
- *name* (*str*) – new entry identifier, must not already exist in embedded files.
- *filename* (*str*) – optional filename or `None`, documentation only, will be set to *name* if `None` or omitted.
- *desc* (*str*) – optional description or `None`, arbitrary documentation text, will be set to *name* if `None` or omitted.

Return type `int`

Returns the index given to the new entry. In the current (April 11, 2017) MuPDF version, this is not reliably true (for this reason we have decided to restrict `embeddedFileDel()` to entries identified by name). Use character string look up to find your entry again. For any error condition, an exception is raised.

`close()`

Release objects and space allocations associated with the document. If created from a file, also closes *filename* (releasing control to the OS).

`outline`

Contains the first *Outline* entry of the document (or `None`). Can be used as a starting point to walk through all outline items. Accessing this property for encrypted, not authenticated documents will raise an `AttributeError`.

Type *Outline*

`isClosed`

`False` / `0` if document is still open, `True` / `1` otherwise. If closed, most other attributes and methods will have been deleted / disabled. In addition, *Page* objects referring to this document (i.e. created with `Document.loadPage()`) and their dependent objects will no longer be usable. For reference purposes, `Document.name` still exists and will contain the filename of the original document (if applicable).

Type `bool`

needsPass

Contains an indicator showing whether the document is encrypted (**True** (1)) or not (**False** (0)). This indicator remains unchanged - even after the document has been authenticated. Precludes incremental saves if set.

Type bool

isEncrypted

This indicator initially equals **needsPass**. After successful authentication, it is set to **False** to reflect the situation.

Type bool

permissions

Shows the permissions to access the document. Contains a dictionary likes this:

```
>>> doc.permissions
{'print': True, 'edit': True, 'note': True, 'copy': True}
```

The keys have the obvious meaning of permissions to print, change, annotate and copy the document, respectively.

Type dict

metadata

Contains the document's meta data as a Python dictionary or **None** (if **isEncrypted** = **True** and **needPass**=**True**). Keys are **format**, **encryption**, **title**, **author**, **subject**, **keywords**, **creator**, **producer**, **creationDate**, **modDate**. All item values are strings or **None**.

Except **format** and **encryption**, the key names correspond in an obvious way to the PDF keys **/Creator**, **/Producer**, **/CreationDate**, **/ModDate**, **/Title**, **/Author**, **/Subject**, and **/Keywords** respectively.

- **format** contains the PDF version (e.g. 'PDF-1.6').
- **encryption** either contains **None** (no encryption), or a string naming an encryption method (e.g. 'Standard V4 R4 128-bit RC4'). Note that an encryption method may be specified **even if needsPass** = **False**. In such cases not all permissions will probably have been granted. Check dictionary **permissions** for details.
- If the date fields contain valid data (which need not be the case at all!), they are strings in the PDF-specific timestamp format "D:<TS><TZ>", where
 - <TS> is the 12 character ISO timestamp **YYYYMMDDhhmmss** (**YYYY** - year, **MM** - month, **DD** - day, **hh** - hour, **mm** - minute, **ss** - second), and
 - <TZ> is a time zone value (time intervall relative to GMT) containing a sign ('+' or '-'), the hour (**hh**), and the minute ('**mm**', note the apostrophies!).
- A Paraguayan value might hence look like **D:20150415131602-04'00'**, which corresponds to the timestamp April 15, 2015, at 1:16:02 pm local time Asuncion.

Type dict

name

Contains the **filename** or **filetype** value with which **Document** was created.

Type str

pageCount

Contains the number of pages of the document. May return 0 for documents with no pages. Function **len(doc)** will also deliver this result.

Type int

openErrCode

If `openErrCode > 0`, errors have occurred while opening / parsing the document, which usually means document structure issues. In this case incremental save cannot be used.

Type int

openErrMsg

Contains either an empty string or the last open error message if `openErrCode > 0`. Together with any other error messages of MuPDF's C library, it will also appear on `SYSERR`.

Type str

embeddedFileCount

Contains the number of files in the embedded / portfolio files list (also known as collection or attached files). If the document is not a PDF, `-1` will be returned.

Type int

Note: For methods that change the structure of a PDF (`insertPDF()`, `select()`, `copyPage()`, `deletePage()` and others), be aware that objects or properties in your program may have been invalidated or orphaned. Examples are *Page* objects and their children (links and annotations), variables holding old page counts, tables of content and the like. Remember to keep such variables up to date or delete orphaned objects.

4.3.1 Remarks on `select()`

Page numbers in the list need not be unique nor be in any particular sequence. This makes the method a versatile utility to e.g. select only the even or the odd pages, re-arrange a document from back to front, duplicate it, and so forth. In combination with text search or extraction you can also omit / include pages with no text or containing a certain text, etc.

You can execute several selections in a row. The document structure will be updated after each method execution.

Any of those changes will become permanent only with a `doc.save()`. If you have de-selected many pages, consider specifying the `garbage` option to eventually reduce the resulting document's size (when saving to a new file).

Also note, that this method **preserves all links, annotations and bookmarks** that are still valid. In other words: deleting pages only deletes references which point to de-selected pages. Page number of bookmarks (outline items) are automatically updated when a TOC is retrieved again with `getToC()`. If a bookmark's destination page happened to be deleted, then its page number in `getToC()` will be set to `-1`.

The results of this method can of course also be achieved using combinations of methods `copyPage()`, `deletePage()` and `movePage()`. While there are many cases, when these methods are more practical, `select()` is easier and safer to use when many pages are involved.

4.3.2 `select()` Examples

In general, any list of integers within the document's page range can be used. Here are some illustrations.

Delete pages with no text:

```
import fitz
doc = fitz.open("any.pdf")
r = list(range(len(doc)))           # list of page numbers

for page in doc:
    if not page.getText():          # page contains no text
```

```
        r.remove(page.number)                # remove page number from list

if len(r) < len(doc):                        # did we actually delete anything?
    doc.select(r)                            # apply the list
doc.save("out.pdf", garbage = 4)             # save result to new PDF, OR

# update the original document ... *** VERY FAST! ***
doc.saveIncr()
```

Create a sub document with only the odd pages:

```
import fitz
doc = fitz.open("any.pdf")
r = list(range(0, len(doc), 2))
doc.select(r)                               # apply the list
doc.save("oddpages.pdf", garbage = 4)        # save sub-PDF of the odd pages
```

Concatenate a document with itself:

```
import fitz
doc = fitz.open("any.pdf")
r = list(range(len(doc)))
r += r                                     # turn PDF into a copy of itself
doc.select(r)
doc.save("any+any.pdf")                    # contains doubled <any.pdf>
```

Create document copy in reverse page order (well, don't try with a million pages):

```
import fitz
doc = fitz.open("any.pdf")
r = list(range(len(doc) - 1, -1, -1))
doc.select(r)
doc.save("back-to-front.pdf")
```

4.3.3 setMetadata() Example

Clear metadata information. If you do this out of privacy / data protection concerns, make sure you save the document as a new file with `garbage > 0`. Only then the old `/Info` object will also be physically removed from the file. In this case you may also want to clear any XML metadata inserted by some PDF editors:

```
>>> import fitz
>>> doc=fitz.open("pymupdf.pdf")
>>> doc.metadata                # look at what we currently have
{'producer': 'rst2pdf, reportlab', 'format': 'PDF 1.4', 'encryption': None, 'author':
'Jorj X. McKie', 'modDate': 'D:20160611145816-04'00'', 'keywords': 'PDF, XPS, EPUB, CBZ',
'title': 'The PyMuPDF Documentation', 'creationDate': 'D:20160611145816-04'00'',
'creator': 'sphinx', 'subject': 'PyMuPDF 1.9.1'}
>>> doc.setMetadata({})        # clear all fields
0
>>> doc.metadata                # look again to show what happened
{'producer': 'none', 'format': 'PDF 1.4', 'encryption': None, 'author': 'none',
'modDate': 'none', 'keywords': 'none', 'title': 'none', 'creationDate': 'none',
'creator': 'none', 'subject': 'none'}
>>> doc._delXmlMetadata()      # clear any XML metadata
0
>>> doc.save("anonymous.pdf", garbage = 4)    # save anonymized doc
0
```

4.3.4 setToC() Example

This shows how to modify or add a table of contents. Also have a look at `csv2toc.py` and `toc2csv.py` in the examples directory:

```
>>> import fitz
>>> doc = fitz.open("test.pdf")
>>> toc = doc.getToC()
>>> for t in toc: print(t)                                # show what we have
...
[1, 'The PyMuPDF Documentation', 1]
[2, 'Introduction', 1]
[3, 'Note on the Name fitz', 1]
[3, 'License', 1]
>>> toc[1][1] += " modified by setToC"                    # modify something
>>> doc.setToC(toc)                                       # replace outline tree
3                                                         # number of bookmarks inserted
>>> for t in doc.getToC(): print(t)                       # demonstrate it worked
...
[1, 'The PyMuPDF Documentation', 1]
[2, 'Introduction modified by setToC', 1]                  # <<< this has changed
[3, 'Note on the Name fitz', 1]
[3, 'License', 1]
```

4.3.5 insertPDF() Examples

(1) Concatenate two documents including their TOCs:

```
doc1 = fitz.open("file1.pdf")          # must be a PDF
doc2 = fitz.open("file2.pdf")          # must be a PDF
pages1 = len(doc1)                     # save doc1's page count
toc1 = doc1.getToC(simple = False)     # save TOC 1
toc2 = doc2.getToC(simple = False)     # save TOC 2
doc1.insertPDF(doc2)                   # doc2 at end of doc1
for t in toc2:                          # increase toc2 page numbers
    t[2] += pages1                     # by old len(doc1)
doc1.setToC(toc1 + toc2)               # now result has total TOC
```

Obviously, similar ways can be found in more general situations. Just make sure that hierarchy levels in a row do not increase by more than one. Inserting dummy bookmarks before and after `toc2` segments would heal such cases.

(2) More examples:

```
# insert 5 pages of doc2, where its page 21 becomes page 15 in doc1
doc1.insertPDF(doc2, from_page = 21, to_page = 25, start_at = 15)

# same example, but pages are rotated and copied in reverse order
doc1.insertPDF(doc2, from_page = 25, to_page = 21, start_at = 15, rotate = 90)

# put copied pages in front of doc1
doc1.insertPDF(doc2, from_page = 21, to_page = 25, start_at = 0)
```

4.3.6 Other Examples

Extract all page-referenced images of a PDF into separate PNG files:

```
for i in range(len(doc)):
    imglist = doc.getPageImageList(i)
    for img in imglist:
```

```
xref = img[0]                # xref number
pix = fitz.Pixmap(doc, xref)  # make pixmap from image
if pix.colourspace.n < 4:    # can be saved as PNG
    pix.writePNG("p%s-%s.png" % (i, xref))
else:                        # CMYK: must convert first
    pix0 = fitz.Pixmap(fitz.csRGB, pix)
    pix0.writePNG("p%s-%s.png" % (i, xref))
    pix0 = None              # free Pixmap resources
pix = None                   # free Pixmap resources
```

Rotate all pages of a PDF:

```
for page in doc:
    page.setRotation(90)
```

4.4 Identity

Identity is just a *Matrix* that performs no action, to be used whenever the syntax requires a *Matrix*, but no actual transformation should take place.

Identity is a constant, an “immutable” object. So, all of its matrix properties are read-only and its methods are disabled.

If you need a do-nothing matrix as a starting point, use `fitz.Matrix(1, 1)` or `fitz.Matrix(0)` instead, like so:

```
>>> fitz.Matrix(0).preTranslate(2, 5)
fitz.Matrix(1.0, 0.0, -0.0, 1.0, 2.0, 5.0)
```

4.5 IRect

IRect is a rectangular bounding box similar to *Rect*, except that all corner coordinates are integers. IRect is used to specify an area of pixels, e.g. to receive image data during rendering. Otherwise, many similarities exist, e.g. considerations concerning emptiness and finiteness of rectangles also apply to IRects.

| Attribute / Method | Short Description |
|----------------------------|--|
| <i>IRect.getRect()</i> | return a <i>Rect</i> with same coordinates |
| <i>IRect.getRectArea()</i> | calculate the area of the rectangle |
| <i>IRect.getArea()</i> | calculate the area of the rectangle |
| <i>IRect.intersect()</i> | common part with another rectangle |
| <i>IRect.translate()</i> | shift rectangle |
| <i>IRect.contains()</i> | checks containment of another object |
| <i>IRect.intersects()</i> | checks for non-empty intersection |
| <i>IRect.normalize()</i> | makes a rectangle finite |
| <i>IRect.height</i> | height of the rectangle |
| <i>IRect.width</i> | width of the rectangle |
| IRect.rect | equals result of method <code>getRect()</code> |
| <i>IRect.top_left</i> | top left point |
| <i>IRect.top_right</i> | top right point |
| <i>IRect.bottom_left</i> | bottom left point |
| <i>IRect.bottom_right</i> | bottom right point |
| <i>IRect.x0</i> | X-coordinate of the top left corner |
| <i>IRect.y0</i> | Y-coordinate of the top left corner |
| <i>IRect.x1</i> | X-coordinate of the bottom right corner |
| <i>IRect.y1</i> | Y-coordinate of the bottom right corner |
| <i>IRect.isInfinite</i> | True if rectangle is infinite |
| <i>IRect.isEmpty</i> | True if rectangle is empty |

Class API

class `IRect`

`__init__(self)`

`__init__(self, x0, y0, x1, y1)`

`__init__(self, irect)`

`__init__(self, list)`

Overloaded constructors. Also see examples below and those for the *Rect* class.

If another *irect* is specified, a **new copy** will be made.

If *list* is specified, it must be a Python sequence type of 4 integers. Non-integer numbers will be truncated, non-numeric entries will be replaced with -1.

The other parameters mean integer coordinates.

`getRect()`

A convenience function returning a *Rect* with the same coordinates as floating point values.

Return type *Rect*

`getRectArea([unit])`

or

`getArea([unit])`

Calculates the area of the rectangle and, with no parameter, equals `abs(IRect)`. Like an empty rectangle, the area of an infinite rectangle is also zero.

Parameters *unit* (*str*) – Specify required unit: respective squares of *px* (pixels, default), *in* (inches), *cm* (centimeters), or *mm* (millimeters).

Return type float

`intersect(ir)`

The intersection (common rectangular area) of the current rectangle and *ir* is calculated and replaces the current rectangle. If either rectangle is empty, the result is also empty. If one of

the rectangles is infinite, the other one is taken as the result - and hence also infinite if both rectangles were infinite.

Parameters *ir* (*IRect*) – Second rectangle.

`translate(tx, ty)`

Modifies the rectangle to perform a shift in x and / or y direction.

Parameters

- *tx* (*int*) – Number of pixels to shift horizontally. Negative values mean shifting left.
- *ty* (*int*) – Number of pixels to shift vertically. Negative values mean shifting down.

`contains(x)`

Checks whether *x* is contained in the rectangle. It may be an *IRect*, *Rect*, “Point” or number. If *x* is an empty rectangle, this is always **True**. Conversely, if the rectangle is empty this is always **False**, if *x* is not an empty rectangle and not a number. If *x* is a number, it will be checked to be one of the four components. *x in rect* and *rect.contains(x)* are equivalent.

Parameters *x* (*IRect* or *Rect* or *Point* or *int*) – the object to check.

Return type *bool*

`intersects(r)`

Checks whether the rectangle and *r* (*IRect* or *Rect*) have a non-empty rectangle in common. This will always be **False** if either is infinite or empty.

Parameters *r* (*IRect* or *Rect*) – the rectangle to check.

Return type *bool*

`normalize()`

Makes sure the rectangle is finite. This is done by shuffling the rectangle corners. After completion of this method, the bottom right corner will indeed be south-eastern to the top left one. See *Rect* for a more detailed discussion on rectangle properties.

`top_left`

Equals *Point(x0, y0)*.

Type *Point*

`top_right`

Equals *Point(x1, y0)*.

Type *Point*

`bottom_left`

Equals *Point(x0, y1)*.

Type *Point*

`bottom_right`

Equals *Point(x1, y1)*.

Type *Point*

`width`

Contains the width of the bounding box. Equals *x1 - x0*.

Type *int*

`height`

Contains the height of the bounding box. Equals *y1 - y0*.

Type *int*

`x0`

X-coordinate of the left corners.


```

    Type int
y0
    Y-coordinate of the top corners.
    Type int
x1
    X-coordinate of the right corners.
    Type int
y1
    Y-coordinate of the bottom corners.
    Type int
isInfinite
    True if rectangle is infinite, False otherwise.
    Type bool
isEmpty
    True if rectangle is empty, False otherwise.
    Type bool

```

4.5.1 Remark

A rectangle's coordinates can also be accessed via index, e.g. `r.x0 == r[0]`, and the `tuple()` and `list()` functions yield sequence objects of its components.

4.5.2 IRect Algebra

The following arithmetic operators have been defined for `IRect` objects (denoted as `ir` in the following). Note that in most binary operations, the second operand may also be of type *Rect*, *Point*, sequences or numbers.

Binary Operators

- **Addition:** `ir + x` where `ir` is an `IRect` and `x` is a number, list / tuple, `Rect` or `IRect`. The result is a new `IRect` with added components of the operands. If `x` is a number, it is added to all components of `ir`.
- **Subtraction:** analogous to addition.
- **Inclusion “|”:** `ir | x` is the new `IRect` that also includes `x` (may be a `Rect`, `IRect` or `Point`).
- **Intersection “&”:** `ir & x` is a new `IRect` containing the area common to `ir` and `x` (`Rect` or `IRect`).
- **Multiplication:** `ir * m` is the new `IRect` resulting from `(ir.rect).transform(m)` for a **matrix** `m` or from multiplication with **number** `m` (coordinate-wise). If the number `m` was not an integer, method `round()` will be applied.
- **Containment Test:** `if x in ir:` tests whether `ir` contains `x` (a number, `Rect` or `IRect`). For a `Rect` or `IRect` this tests whether its area is contained in `ir`. If `x` is a number it tests whether `x` is one of the 4 coordinates.
- **Comparison:** `ir1 == ir2` is `True` if the tuples of coordinates are equal (not only if they are the same object!). However, `rect == irect` is always `False` because of the different object types.

Unary Operators

- `-ir` is a new copy of `ir` with negated components.
- `+ir` is a new copy of `ir`.
- `bool(ir)` is `False` for `IRect(0, 0, 0, 0)` and `True` otherwise.
- `abs(ir)` is equal to `ir.getArea()`.

4.5.3 Examples

Algebra provides handy ways to perform inclusion and intersection checks between `Rects`, `IRects` and `Points`.

Example 1:

```
>>> ir = fitz.IRect(10, 10, 410, 610)
>>> ir
fitz.IRect(10, 10, 410, 610)
>>> ir.height
600
>>> ir.width
400
>>> ir.getArea('mm')      # calculate area in square millimeters
29868.51852
```

Example 2:

```
>>> m = fitz.Matrix(45)
>>> ir = fitz.IRect(10, 10, 410, 610)
>>> ir * m                  # rotate rectangle by 45 degrees
fitz.IRect(-425, 14, 283, 722)
>>>
>>> ir | fitz.Point(5, 5)   # enlarge rectangle to contain a point
fitz.IRect(5, 5, 410, 610)
>>>
>>> ir + 5                  # shift the rect by 5 points
fitz.IRect(15, 15, 415, 615)
>>>
>>> ir & fitz.Rect(0.0, 0.0, 15.0, 15.0)
fitz.IRect(10, 10, 15, 15)
```

Example 3:

```
>>> # test whether two rectangle are disjoint
>>> if not r1.intersects(r2): print("disjoint rectangles")
>>>
>>> # test whether r2 contains x (x is Point, Rect, IRect or number)
>>> if r2.contains(x): print("x is contained in r2")
>>>
>>> # or even simpler:
>>> if x in r2: print("x is contained in r2")
```

4.6 Link

Represents a pointer to somewhere (this document, other documents, the internet). Links exist per document page, and they are forward-chained to each other, starting from an initial link which is accessible by the `Page.loadLinks()` method.

There is a parent-child relationship between a link and its page. If the page object becomes unusable (closed document, any document structure change, etc.), then so does every of its existing link objects - an exception is raised saying that the object is “orphaned”, whenever a link property or method is accessed.

| Attribute | Short Description |
|------------------------|--|
| <i>Link.rect</i> | clickable area in untransformed coordinates. |
| <i>Link.uri</i> | link destination |
| <i>Link.isExternal</i> | link destination |
| <i>Link.next</i> | points to next link |
| <i>Link.dest</i> | points to link destination details |

Class API

class Link

rect

The area that can be clicked in untransformed coordinates.

Type *Rect*

isExternal

A bool specifying whether the link target is outside (**True**) of the current document.

Type bool

uri

A string specifying the link target. The meaning of this property should be evaluated in conjunction with property **isExternal**. The value may be **None**, in which case **isExternal** == **False**. If **uri** starts with **file://**, **mailto:**, or an internet resource name, **isExternal** is **True**. In all other cases **isExternal** == **False** and **uri** points to an internal location. In case of PDF documents, this should either be **#nnnn** to indicate a 1-based (!) page number **nnnn**, or a named location. The format varies for other document types, e.g. **uri** = **'../FixedDoc.fdoc#PG_2_LNK_1'** for page number 2 (1-based) in an XPS document.

Type str

next

The next Link or **None**

Type Link

dest

The link destination details object.

Type *linkDest*

4.7 linkDest

Class representing the *dest* property of an outline entry or a link. Describes the destination to which such entries point.

| Attribute | Short Description |
|---------------------------|-------------------------------------|
| <i>linkDest.dest</i> | destination |
| <i>linkDest.fileSpec</i> | file specification (path, filename) |
| <i>linkDest.flags</i> | descriptive flags |
| <i>linkDest.isMap</i> | is this a MAP? |
| <i>linkDest.isUri</i> | is this a URI? |
| <i>linkDest.kind</i> | kind of destination |
| <i>linkDest.lt</i> | top left coordinates |
| <i>linkDest.named</i> | name if named destination |
| <i>linkDest.newWindow</i> | name of new window |
| <i>linkDest.page</i> | page number |
| <i>linkDest.rb</i> | bottom right coordinates |
| <i>linkDest.uri</i> | URI |

Class API

class linkDest

dest

Target destination name if *linkDest.kind* is *LINK_GOTOR* and *linkDest.page* is -1.

Type str

fileSpec

Contains the filename and path this link points to, if *linkDest.kind* is *LINK_GOTOR* or *LINK_LAUNCH*.

Type str

flags

A bitfield describing the validity and meaning of the different aspects of the destination. As far as possible, link destinations are constructed such that e.g. *linkDest.lt* and *linkDest.rb* can be treated as defining a bounding box. But the flags indicate which of the values were actually specified, see *Link Destination Flags*.

Type int

isMap

This flag specifies whether to track the mouse position when the URI is resolved. Default value: False.

Type bool

isUri

Specifies whether this destination is an internet resource (as opposed to e.g. a local file specification in URI format).

Type bool

kind

Indicates the type of this destination, like a place in this document, a URI, a file launch, an action or a place in another file. Look at *Enumerations* to see the names and numerical values.

Type int

lt

The top left *Point* of the destination.

Type *Point*

named

This destination refers to some named action to perform (e.g. a javascript, see *Adobe PDF Reference 1.7*). Standard actions provided are *NextPage*, *PrevPage*, *FirstPage*, and *LastPage*.

Type str

newWindow

If true, the destination should be launched in a new window.

Type bool

page

The page number (in this or the target document) this destination points to. Only set if *linkDest.kind* is *LINK_GOTOR* or *LINK_GOTO*. May be -1 if *linkDest.kind* is *LINK_GOTOR*. In this case *linkDest.dest* contains the **name** of a destination in the target document.

Type int

rb

The bottom right *Point* of this destination.

Type *Point*

uri

The name of the URI this destination points to.

Type str

4.8 Matrix

Matrix is a row-major 3x3 matrix used by image transformations in MuPDF (which complies with the respective concepts laid down in the *Adobe PDF Reference 1.7*). With matrices you can manipulate the rendered image of a page in a variety of ways: (parts of) the page can be rotated, zoomed, flipped, sheared and shifted by setting some or all of just six float values.

Since all points or pixels live in a two-dimensional space, one column vector of that matrix is a constant unit vector, and only the remaining six elements are used for manipulations. These six elements are usually represented by [a, b, c, d, e, f]. Here is how they are positioned in the matrix:

$$\begin{bmatrix} a & b & 0 \\ c & d & 0 \\ e & f & 1 \end{bmatrix}$$

Please note:

- the below methods are just convenience functions - everything they do, can also be achieved by directly manipulating the six numerical values
- all manipulations can be combined - you can construct a matrix that rotates **and** shears **and** scales **and** shifts, etc. in one go. If you however choose to do this, do have a look at the **remarks** further down or at the *Adobe PDF Reference 1.7*.

| Method / Attribute | Description |
|------------------------------|----------------------------------|
| <i>Matrix.preRotate()</i> | perform a rotation |
| <i>Matrix.preScale()</i> | perform a scaling |
| <i>Matrix.preShear()</i> | perform a shearing (skewing) |
| <i>Matrix.preTranslate()</i> | perform a translation (shifting) |
| <i>Matrix.concat()</i> | perform a matrix multiplication |
| <i>Matrix.invert()</i> | calculate the inverted matrix |
| <i>Matrix.a</i> | zoom factor X direction |
| <i>Matrix.b</i> | shearing effect Y direction |
| <i>Matrix.c</i> | shearing effect X direction |
| <i>Matrix.d</i> | zoom factor Y direction |
| <i>Matrix.e</i> | horizontal shift |
| <i>Matrix.f</i> | vertical shift |

Class API

class Matrix

```
__init__(self)
__init__(self, zoom-x, zoom-y)
__init__(self, shear-x, shear-y, 1)
__init__(self, a, b, c, d, e, f)
__init__(self, matrix)
__init__(self, degree)
__init__(self, list)
```

Overloaded constructors.

Without parameters, `Matrix(0.0, 0.0, 0.0, 0.0, 0.0, 0.0)` will be created.

`zoom-*` and `shear-*` specify zoom or shear values (float), respectively.

`matrix` specifies another `Matrix` from which a **new copy** will be made.

Float value `degree` specifies the creation of a rotation matrix.

Python sequence `list` (list, tuple, etc.) must contain exactly 6 values when specified. Other than that **no error checking is done**, and invalid entries will receive a value of `-1.0`.

`fitz.Matrix(1, 1)`, `fitz.Matrix(0.0)` and `fitz.Matrix(fitz.Identity)` create modifiable versions of the *Identity* matrix, which looks like `[1, 0, 0, 1, 0, 0]`.

`preRotate(deg)`

Modify the matrix to perform a counter-clockwise rotation for positive `deg` degrees, else clockwise. The matrix elements of an identity matrix will change in the following way:

`[1, 0, 0, 1, 0, 0] -> [cos(deg), sin(deg), -sin(deg), cos(deg), 0, 0]`.

Parameters `deg (float)` – The rotation angle in degrees (use conventional notation based on $\text{Pi} = 180$ degrees).

`preScale(sx, sy)`

Modify the matrix to scale by the zoom factors `sx` and `sy`. Has effects on attributes `a` thru `d` only: `[a, b, c, d, e, f] -> [a*sx, b*sx, c*sy, d*sy, e, f]`.

Parameters

- `sx (float)` – Zoom factor in X direction. For the effect see description of attribute `a`.

- `sy (float)` – Zoom factor in Y direction. For the effect see description of attribute `d`.

`preShear(sx, sy)`

Modify the matrix to perform a shearing, i.e. transformation of rectangles into parallelograms (rhomboids). Has effects on attributes `a` thru `d` only: `[a, b, c, d, e, f] -> [c*sy, d*sy, a*sx, b*sx, e, f]`.

Parameters

- `sx (float)` – Shearing effect in X direction. See attribute `c`.
- `sy (float)` – Shearing effect in Y direction. See attribute `b`.

`preTranslate(tx, ty)`

Modify the matrix to perform a shifting / translation operation along the x and / or y axis. Has effects on attributes `e` and `f` only: `[a, b, c, d, e, f] -> [a, b, c, d, tx*a + ty*c, tx*b + ty*d]`.

Parameters

- `tx (float)` – Translation effect in X direction. See attribute `e`.
- `ty (float)` – Translation effect in Y direction. See attribute `f`.

`concat(m1, m2)`

Calculate the matrix product `m1 * m2` and store the result in the current matrix. Any of `m1` or `m2` may be the current matrix. Be aware that matrix multiplication is not commutative. So the sequence of `m1`, `m2` is important.

Parameters

- `m1 (Matrix)` – First (left) matrix.
- `m2 (Matrix)` – Second (right) matrix.

`invert(m)`

Calculate the matrix inverse of `m` and store the result in the current matrix. Returns 1 if `m` is not invertible (“degenerate”). In this case the current matrix **will not change**. Returns 0 if `m` is invertible, and the current matrix is replaced with the inverted `m`.

Parameters `m (Matrix)` – Matrix to be inverted.

Return type `int`

a

Scaling in X-direction (**width**). For example, a value of 0.5 performs a shrink of the **width** by a factor of 2. If `a < 0`, a left-right flip will (additionally) occur.

Type `float`

b

Causes a shearing effect: each `Point(x, y)` will become `Point(x, y - b*x)`. Therefore, looking from left to right, e.g. horizontal lines will be “tilt” - downwards if `b > 0`, upwards otherwise (`b` is the tangens of the tilting angle).

Type `float`

c

Causes a shearing effect: each `Point(x, y)` will become `Point(x - c*y, y)`. Therefore, looking upwards, vertical lines will be “tilt” - to the left if `c > 0`, to the right otherwise (`c` is the tangens of the tilting angle).

Type `float`

d

Scaling in Y-direction (**height**). For example, a value of 1.5 performs a stretch of the **height** by 50%. If `d < 0`, an up-down flip will (additionally) occur.

Type `float`

e
Causes a horizontal shift effect: Each `Point(x, y)` will become `Point(x + e, y)`. Positive (negative) values of **e** will shift right (left).

Type float

f
Causes a vertical shift effect: Each `Point(x, y)` will become `Point(x, y - f)`. Positive (negative) values of **f** will shift down (up).

Type float

4.8.1 Remarks 1

For a matrix `m`, properties **a** to **f** can also be accessed by index, e.g. `m.a == m[0]` and `m[0] = 1` has the same effect as `m.a = 1`. The `tuple()` and `list()` functions yield sequence objects of its components.

Language constructs like `x in m` will check whether number `x` is one of the six components.

4.8.2 Remarks 2

Changes of matrix properties and execution of matrix methods can be executed consecutively. This is the same as multiplying the respective matrices.

Matrix multiplications are **not commutative** - changing the execution sequence in general changes the result. So it can quickly become unclear which result a transformation will yield.

In order to keep results foreseeable for a series of transformations, Adobe recommends the following approach (*Adobe PDF Reference 1.7*, page 206):

1. Shift (“translate”)
2. Rotate
3. Scale or shear (“skew”)

4.8.3 Matrix Algebra

The following arithmetic operators have been defined for the `Matrix` class. In what follows, `m`, `m1`, `m2` denote matrices:

- **Binary Operators**

- **Addition:** `m1 + m2` is the new matrix containing `[m1.a + m2.a, ..., m1.f + m2.f]`. `m2` may also be any sequence of 6 floats.
- **Subtraction:** analogous to addition.
- **Multiplication:** `m1 * m2` is a new matrix calculated as `concat(m1, m2)`. If `m2` is a number however, it will simply multiply the coordinates.
- **Comparison:** `m1 == m2` is `True` if the coordinates are equal (not only if they are the same object!). This tests equality between floats, so be wary of artifact differences.

- **Unary Operators:**

- `-m` is a copy of `m` with negated components.
- `+m` is a copy of `m`.
- `~m` is the inverted matrix of `m`, i.e. `m * ~m = ~m * m = fitz.Identity`. Due to rounding issues, this equality will mostly not be exact, however. See examples below. If `m` is degenerate (not invertible), `~m` will be `fitz.Matrix(0, 0, 0, 0, 0, 0)` or `bool(~m) = False`.

- **Absolute Value:** `abs(m)` is a float containing the Euclidean norm of `m`. Typically used for testing whether two matrices are “almost equal”, like `abs(m1 - m2) < epsilon`.
- **Zero Test:** `bool(m)` is `False` for `Matrix(0, 0, 0, 0, 0, 0)` and `True` otherwise. Can be used to test invertibility of matrices: `not bool(~m)` means `m` is degenerate.

This makes the following operations possible:

```
>>> m45p = fitz.Matrix(45)           # rotate 45 degrees clockwise
>>> m45m = fitz.Matrix(-45)          # rotate 45 degrees counterclockwise
>>> m90p = fitz.Matrix(90)           # rotate 90 degrees clockwise
>>>
>>> abs(m45p * ~m45p - fitz.Identity) # should be (close to) zero:
8.429369702178807e-08
>>>
>>> abs(m90p - m45p * m45p)          # should be (close to) zero:
8.429369702178807e-08
>>>
>>> abs(m45p * m45m - fitz.Identity) # should be (close to) zero:
2.1073424255447017e-07
>>>
>>> abs(m45p - ~m45m)                # should be (close to) zero:
2.384185791015625e-07
>>>
>>> m90p * m90p * m90p * m90p        # should be 360 degrees = fitz.Identity
fitz.Matrix(1.0, -0.0, 0.0, 1.0, 0.0, 0.0)
```

4.8.4 Examples

Here are examples to illustrate some of the effects achievable. The following pictures start with a page of the PDF version of this help file. We show what happens when a matrix is being applied (though always full pages are created, only parts are displayed here to save space).

This is the original page image:

Classes

Matrix

Matrix is a row-major 3x3 matrix used for representing transformations of coordinates throughout MuPDF.

Since all points or pixels reside in a two-dimensional space, one column vector of the matrix is the constant unit vector, and only the remaining six elements may vary. These six elements are usually represented by `[a, b, c, d, e, f]`. Here is how they are positioned in the matrix:

$$\begin{bmatrix} a & b & 0 \\ c & d & 0 \\ e & f & 1 \end{bmatrix}$$

It should be noted, that the below methods are just convenience functions. Each of them manipulates some of the six matrix elements in a specific way. By directly changing `[a, b, c, d, e, f]`, any of these functions can be replaced.

4.8.5 Shifting

We transform it with a matrix where `e = 100` (right shift by 100 pixels).

Classes

Matrix is a row-major 3x3 matrix used for representing transformations of coordinates throughout MuPDF. Since all points or pixels reside in a two-dimensional space, one column vector of the matrix is the constant unit vector, and only the remaining six elements may vary. These six elements are usually represented by $[a, b, c, d, e, f]$. Here is how they are positioned in the matrix:

$$\begin{bmatrix} a & b & 0 \\ c & d & 0 \\ e & f & 1 \end{bmatrix}$$

Next we do a down shift by 100 pixels: $f = 100$.

Classes

Matrix

Matrix is a row-major 3x3 matrix used for representing transformations of coordinates throughout MuPDF.

Since all points or pixels reside in a two-dimensional space, one column vector of the matrix is the constant unit vector, and only the remaining six elements may vary. These six elements are usually represented by $[a, b, c, d, e, f]$. Here is how they are positioned in the matrix:

$$\begin{bmatrix} a & b & 0 \\ c & d & 0 \\ e & f & 1 \end{bmatrix}$$

4.8.6 Flipping

Flip the page left-right ($a = -1$).

Classes

Matrix

Matrix is a row-major 3x3 matrix used for representing transformations of coordinates throughout MuPDF. Since all points or pixels reside in a two-dimensional space, one column vector of the matrix is the constant unit vector, and only the remaining six elements may vary. These six elements are usually represented by $[a, b, c, d, e, f]$. Here is how they are positioned in the matrix:

$$\begin{bmatrix} 0 & d & a \\ 0 & b & c \\ 1 & f & e \end{bmatrix}$$

Flip up-down ($d = -1$).

$$\begin{bmatrix} a & b & 1 \\ c & d & 0 \\ e & f & 0 \end{bmatrix}$$

$[a, b, c, d, e, f]$. Here is how they are positioned in the matrix:

vector' and only the remaining six elements may vary. These six elements are usually represented by $[a, b, c, d, e, f]$ since all points or pixels reside in a two-dimensional space; one column vector of the matrix is the constant unit vector.

Matrix is a row-major 3x3 matrix used for representing transformations of coordinates throughout MuPDF.

Matrix

Classes

4.8.7 Shearing

First a shear in Y direction ($b = 0.5$).

Classes

Matrix

Matrix is a row-major 3x3 matrix used image transformations in MuPDF. With matrices you can manipulate the rendered image of a page in a variety of ways: (parts of) pages can be rotated, zoomed, flipped, sheared and shifted by setting some or all of just six numerical values.

Since all points or pixels live in a two-dimensional space, one column vector of that matrix is a constant unit vector, and only the remaining six elements are used for manipulations. These six elements are usually represented by $[a, b, c, d, e, f]$. Here is how they are positioned in the matrix:

$$\begin{bmatrix} a & b & 0 \\ c & d & 0 \\ e & f & 1 \end{bmatrix}$$

It should be noted, that

- the below methods are just convenience functions. Even manipulating $[a, b, c, d, e, f]$
- all manipulations can be combined - you can combine

Methods

Matrix, ...

Second a shear in X direction ($c = 0.5$).

Classes

Matrix

Matrix is a row-major 3x3 matrix used for image transformations in MuPDF. With matrices you can manipulate the rendered image of a page in a variety of ways: (parts of) pages can be rotated, zoomed, flipped, sheared and shifted by setting some or all of just six numerical values.

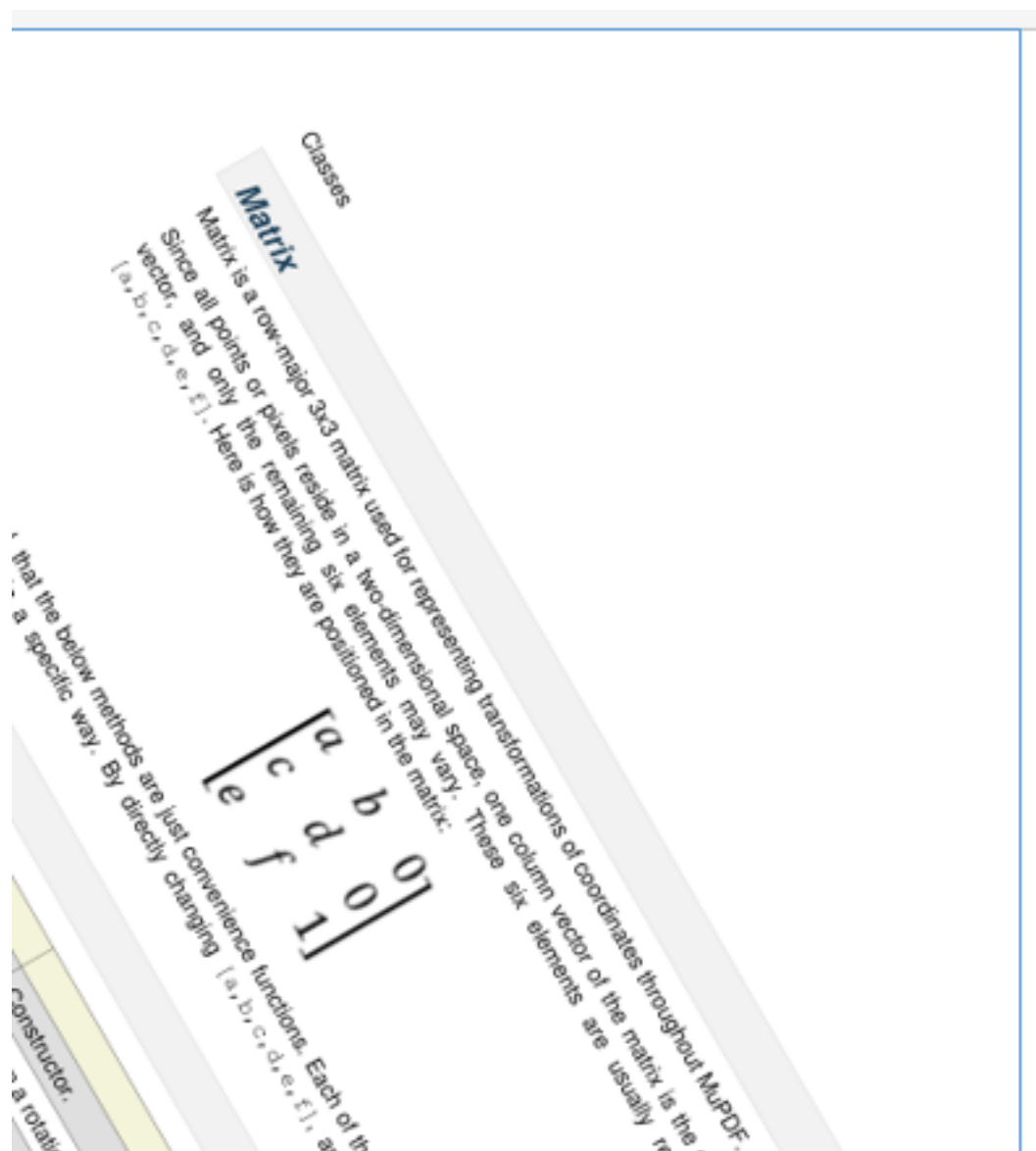
Since all points or pixels live in a two-dimensional space, one column vector of that matrix is a constant unit vector, and only the remaining six elements are used for manipulations. These six elements are usually represented by `[a, b, c, d, e, f]`. Here is how they are positioned in the matrix:

$$\begin{bmatrix} a & b & 0 \\ c & d & 0 \\ e & f & 1 \end{bmatrix}$$

It should be noted, that

4.8.8 Rotating

Finally a rotation by 30 clockwise degrees (`preRotate(-30)`).



4.9 Outline

`outline` (or “bookmark”), is a property of `Document`. If not `None`, it stands for the first outline item of the document. Its properties in turn define the characteristics of this item and also point to other outline items in “horizontal” or downward direction. The full tree of all outline items for e.g. a conventional table of contents (TOC) can be recovered by following these “pointers”.

| Method / Attribute | Short Description |
|--|--|
| <i><code>Outline.saveText()</code></i> | prints a conventional TOC to a file |
| <i><code>Outline.saveXML()</code></i> | prints an XML-like TOC to a file |
| <i><code>Outline.down</code></i> | next item downwards |
| <i><code>Outline.next</code></i> | next item same level |
| <i><code>Outline.page</code></i> | page number (0-based) |
| <i><code>Outline.title</code></i> | title |
| <i><code>Outline.uri</code></i> | string further specifying the outline target |
| <i><code>Outline.isExternal</code></i> | target is outside this document |
| <i><code>Outline.is_open</code></i> | whether sub-outlines are open or collapsed |
| <code>Outline.isOpen</code> | whether sub-outlines are open or collapsed |
| <i><code>Outline.dest</code></i> | points to link destination details |

Class API

class `Outline`

`down`

The next outline item on the next level down. Is `None` if the item has no kids.

Type *`Outline`*

`next`

The next outline item at the same level as this item. Is `None` if this is the last one in its level.

Type *`Outline`*

`page`

The page number (0-based) this bookmark points to.

Type `int`

`title`

The item’s title as a string or `None`.

Type `str`

`is_open`

Or `isOpen` - an indicator showing whether any sub-outlines should be expanded (`True`) or be collapsed (`False`). This information should be interpreted by PDF display software accordingly.

Type `bool`

`saveText(filename)`

The chain of outline items is being processed and printed to the file `filename` as a conventional table of contents. Each line of this file has the format `<tab>...<tab><title><tab><page#>`, where the number of leading tabs is (n-1), with n equal to the outline hierarchy level of the entry. Page numbers are 1-based in this case. `page = -1` can occur if the destination is outside this document or undefined (`uri == None`).

Parameters `filename` (*`str`*) – Name of the file to write to.

`saveXML(filename)`

The chain of outline items is being processed and printed to a file `filename` as an XML-like

table of contents. Each line of this file has the format `<outline title="..." page="n"/>`, if the entry has no children. Otherwise the format is `<outline title="..." page="n">`, and child entries will follow. The parent entry will be finished by a line containing `</outline>`.

Parameters `filename` (*str*) – Name of the file to write to.

isExternal

A bool specifying whether the target is outside (`True`) of the current document.

Type `bool`

uri

A string specifying the link target. The meaning of this property should be evaluated in conjunction with `isExternal`. The value may be `None`, in which case `isExternal == False`. If `uri` starts with `file://`, `mailto:`, or an internet resource name, `isExternal` is `True`. In all other cases `isExternal == False` and `uri` points to an internal location. In case of PDF documents, this should either be `#nnnn` to indicate a 1-based (!) page number `nnnn`, or a named location. The format varies for other document types, e.g. `uri = '../FixedDoc.fdoc#PG_21_LNK_84'` for page number 21 (1-based) in an XPS document.

Type `str`

dest

The link destination details object.

Type *linkDest*

4.10 Page

Class representing a document page. A page object is created by *Document.loadPage()* or, equivalently, via indexing the document like `doc[n]` - it has no independent constructor.

There is a parent-child relationship between a document and its pages. If the document is closed or deleted, all page objects (and their respective children, too) in existence will become unusable. If a page property or method is being used, an exception is raised saying that the page object is “orphaned”.

Several page methods have a *Document* counterpart for convenience. At the end of this chapter you will find a synopsis.

| Method / Attribute | Short Description |
|-----------------------------|---------------------------------------|
| <i>Page.bound()</i> | rectangle (mediabox) of the page |
| <i>Page.deleteAnnot()</i> | PDF only: delete an annotation |
| <i>Page.deleteLink()</i> | PDF only: delete a link |
| <i>Page.drawLine()</i> | PDF only: draw a line |
| <i>Page.drawCircle()</i> | PDF only: draw a circle |
| <i>Page.drawBezier()</i> | PDF only: draw a cubic Bézier curve |
| <i>Page.drawCurve()</i> | PDF only: draw a special Bézier curve |
| <i>Page.drawOval()</i> | PDF only: draw an oval / ellipse |
| <i>Page.drawSector()</i> | PDF only: draw a circular sector |
| <i>Page.drawPolyline()</i> | PDF only: draw a polyline |
| <i>Page.drawRect()</i> | PDF only: draw a rectangle |
| <i>Page.getFontList()</i> | PDF only: get list of used fonts |
| <i>Page.getImageList()</i> | PDF only: get list of used images |
| <i>Page.getLinks()</i> | get all links |
| <i>Page.getPixmap()</i> | create a <i>Pixmap</i> |
| <i>Page.getText()</i> | extract the page's text |
| <i>Page.insertImage()</i> | PDF only: insert an image |
| <i>Page.insertLink()</i> | PDF only: insert a new link |
| <i>Page.insertText()</i> | PDF only: insert text |
| <i>Page.insertTextbox()</i> | PDF only: insert a text box |
| <i>Page.searchFor()</i> | search for a string |
| <i>Page.setRotation()</i> | PDF only: set page rotation |
| <i>Page.updateLink()</i> | PDF only: modify a link |
| <i>Page.firstAnnot</i> | first <i>Annot</i> on the page |
| <i>Page.firstLink</i> | first <i>Link</i> on the page |
| <i>Page.number</i> | page number |
| <i>Page.parent</i> | owning document object |
| <i>Page.rect</i> | rectangle (mediabox) of the page |
| <i>Page.rotation</i> | PDF only: page rotation |

Class API

class Page

bound()

Determine the rectangle (“mediabox”, before transformation) of the page.

Return type *Rect*

rect

Contains the rectangle (“mediabox”, before transformation) of the page. Same as result of method `bound()`.

Type *Rect*

deleteAnnot(*annot*)

PDF only: Delete the specified annotation from the page and (for all document types) return the next one.

Parameters *annot* (*Annot*) – the annotation to be deleted.

Return type *Annot*

Returns the next annotation of the deleted one.

deleteLink(*linkdict*)

PDF only: Delete the specified link from the page. The parameter must be a dictionary of format as provided by the `getLinks()` method (see below).

Parameters *linkdict* (*dict*) – the link to be deleted.

`insertLink(linkdict)`

PDF only: Insert a new link on this page. The parameter must be a dictionary of format as provided by the `getLinks()` method (see below).

Parameters `linkdict` (*dict*) – the link to be inserted.

`updateLink(linkdict)`

PDF only: Modify the specified link. The parameter must be a dictionary of format as provided by the `getLinks()` method (see below).

Parameters `linkdict` (*dict*) – the link to be modified.

`getLinks()`

Retrieves **all** links of a page.

Return type list

Returns A list of dictionaries. The entries are in the order as specified during PDF generation. For a description of the dictionary entries see below. Always use this method if you intend to make changes to the links of a page.

`insertText(point, text = text, fontsize = 11, fontname = "Helvetica", fontfile = None, color = (0, 0, 0), wordspacing = 0, rotate = 0, overlay = True)`

PDF only: Insert text.

Parameters

- **point** (*Point*) – the bottom-left position of the first **text** character in pixels. `point.x` specifies the distance from left border, `point.y` the distance from top of page. This is independent from text orientation as requested by **rotate**. However, there must always be sufficient room “above”, which can mean the distance from any of the four page borders.
- **text** (*str or sequence*) – the text to be inserted. May be specified as either a string type or as a sequence type. For sequences, or strings containing line breaks `\n`, several lines will be inserted. No care will be taken if lines are too wide, but the number of inserted lines will be limited by “vertical” space on the page (in the sense of reading direction as established by the **rotate** parameter). Any rest of **text** is discarded - the return code however contains the number of inserted lines. Only single byte character codes are currently supported.
- **wordspacing** (*float*) – amount to be added to the width of each space character `chr(32)` in unscaled units. This can be used to control the space between words. In order to e.g. double the word distances, specify the product of glyph width of `chr(32)` and **fontsize**. Negative numbers are possible. **Caution:** this parameter will be set once per method invocation, and therefore applies to **all lines** that are contained in **text**.
- **rotate** (*int*) – determines whether to rotate the text. Acceptable values are multiples of 90 degrees. Default is 0 (no rotation), meaning horizontal text lines oriented from left to right. 180 means text is shown upside down from **right to left**. 90 means counter-clockwise rotation, text running **upwards**. 270 (or -90) means clockwise rotation, text running **downwards**. In any case, **point** specifies the bottom-left coordinates of the first character’s rectangle. Multiple lines, if present, always follow the reading direction established by this parameter. So line 2 is located **above** line 1 in case of **rotate** = 180, etc.

Return type int

Returns number of lines inserted.

For a description of the other parameters see *Common Parameters*.

`insertTextbox(rect, buffer, fontsize = 11, fontname = "Helvetica", fontfile = None, color = (0, 0, 0), expandtabs = 8, align = TEXT_ALIGN_LEFT, charwidths = None, rotate = 0, overlay = True)`

PDF only: Insert text into the specified rectangle. This is a convenience method that invokes `insertText()` after formatting the text contained in `buffer` as follows. The text will be split into lines and words and then filled into the available space, starting from one of the four rectangle corners, depending on `rotate`. Line feeds will be respected as well as multiple spaces will be.

Parameters

- `rect` (*Rect*) – the area to use. It must be finite, not empty and completely contained in the page.
- `buffer` – the text to be inserted. Must be specified as a string or a sequence of strings. Line breaks are respected also when occurring in a sequence entry.
- `align` (*int*) – align each text line. Default is 0 (left). Centered, right and justified are the other supported options, see *Text Alignment*.
- `expandtabs` (*int*) – controls handling of tab characters `\t` using the `string.expandtabs()` method **per each line**.
- `charwidths` (*sequence*) – specify a float sequence (list, tuple, ...) of character widths for the specified font. If omitted, it will be created by the function on **each invocation**. For efficiency, provide this parameter if you insert several text boxes with the same font. Use low-level function *Document._getCharWidths()* to create one for your font. Only single byte character codes are currently supported. Results are unpredictable, if larger codes occur (and the `charwidths` sequence is longer than 256).
- `rotate` (*int*) – requests text to be rotated in the rectangle. This value must be a multiple of 90 degrees. Default is 0 (no rotation). Effectively, four different values are processed: 0, 90, 180 and 270 (= -90), each causing the text to start in a different rectangle corner. Bottom-left is 90, bottom-right is 180, and -90 / 270 is top-right. See method `insertText()` for more details.

Returns

If positive or zero: successful execution. The value returned is the unused rectangle space in pixels. This may safely be ignored - or be used to optimize the rectangle, position subsequent items, etc.

If negative: no execution. The value returned is the space deficit to store the text. Increase the rectangle, decrease `fontsize`, decrease text length, etc.

Return type float

For a description of the other parameters see *Common Parameters*.

`drawLine(p1, p2, color = (0, 0, 0), width = 1, dashes = None, roundCap = True, overlay = True)`

PDF only: Draw a line.

Parameters

- `p1` (*Point*) – starting point of the line.
- `p2` (*Point*) – end point of the line.

For a description of the other parameters see *Common Parameters*.

`drawCircle(center, radius, color = (0, 0, 0), fill = None, width = 1, dashes = None, roundCap = True, overlay = True)`

PDF only: Draw a circle around `center` with a radius of `radius`.

Parameters

- `center` (*Point*) – center of the circle.
- `radius` (*float*) – radius of the circle.

For a description of the other parameters see [Common Parameters](#).

`drawOval(rect, color = (0, 0, 0), fill = None, width = 1, dashes = None, roundCap = True, overlay = True)`

PDF only: Draw an oval (ellipse) within the given rectangle.

Parameters `rect` ([Rect](#)) – the rectangle in which the oval is to be drawn. The width-to-height proportion determines the oval’s shape, and the oval touches the middle point of each rectangle side. If the rectangle is a square, a regular circle is drawn like with method `drawCircle()`.

For a description of the other parameters see [Common Parameters](#).

`drawSector(center, point, angle, color = (0, 0, 0), fill = None, width = 1, dashes = None, roundCap = True, fullSector = True, overlay = True, closePath = False)`

PDF only: Draw a circular sector, optionally connecting the arc to the circle’s center (like a piece of pie).

Parameters

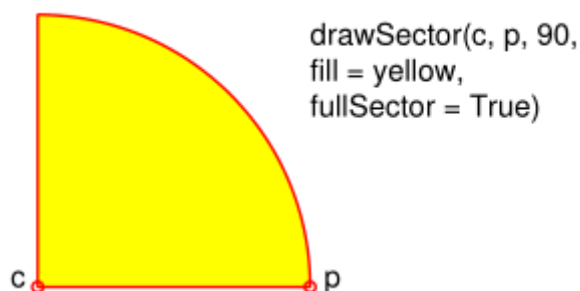
- **center** ([Point](#)) – the center of the circle.
- **point** ([Point](#)) – one of the two end points of the pie’s arc segment. The other one is calculated from the **angle**.
- **angle** (*float*) – the angle of the sector in degrees. Used to calculate the other end point of the arc. Depending on its sign, the arc is drawn counter-clockwise (positive) or clockwise.
- **fullSector** (*bool*) – whether to draw connecting lines from the ends of the arc to the circle center. If a fill color is specified, the full “pie” is colored, otherwise just the sector. If true, parameter `closePath` has no effect.

For a description of the other parameters see [Common Parameters](#).

Returns the other end point of the arc. Can be used as starting point for another invocation to create logically connected pies charts.

Return type [Point](#)

`drawSector()` examples:



```
drawPolyline(points, color = (0, 0, 0), fill = None, width = 1, dashes = None, roundCap =
    True, overlay = True, closePath = False)
```

PDF only: Draw several connected lines defined by a sequence of points.

Parameters *points* (*list or tuple*) – a sequence of *Point* objects defining the edges. The sequence length must be at least two.

For a description of the other parameters see *Common Parameters*.

```
drawBezier(p1, p2, p3, p4, color = (0, 0, 0), fill = None, width = 1, dashes = None, roundCap =
    True, overlay = True, closePath = False)
```

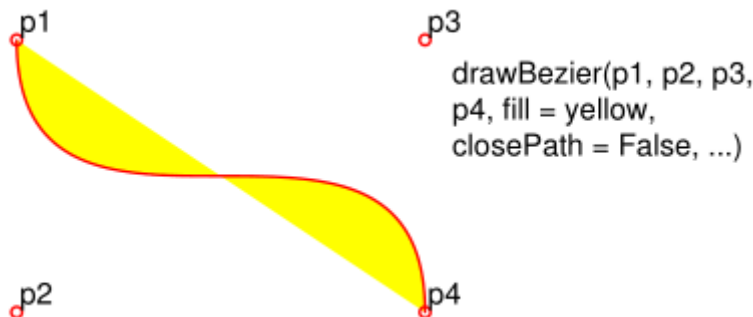
PDF only: Draw a cubic Bézier curve from *p1* to *p4* with the control points *p2* and *p3*.

Parameters

- *p1* (*Point*) – is the starting point of the curve.
- *p2* (*Point*) – control point.
- *p3* (*Point*) – control point.
- *p4* (*Point*) – is the end point of the curve.

For a description of the other parameters see *Common Parameters*.

`drawBezier()` example:



```
drawCurve(p1, p2, p3, color = (0, 0, 0), fill = None, width = 1, dashes = None, roundCap =
    True, overlay = True, closePath = False)
```

PDF only: This is a special case of `drawBezier()`. Draw a cubic Bézier curve from *p1* to *p3*. On each of the lines from *p1* to *p2* and from *p2* to *p3* one control point is generated. This guaranties that the curve's curvature does not change its sign. The control points are defined as $c1 = p1 + (p2 - p1) * k$ and $c2 = p3 + (p2 - p3) * k$, where $k = 0.552$, a value giving excellent approximations for quarter circle / ellipse arcs.

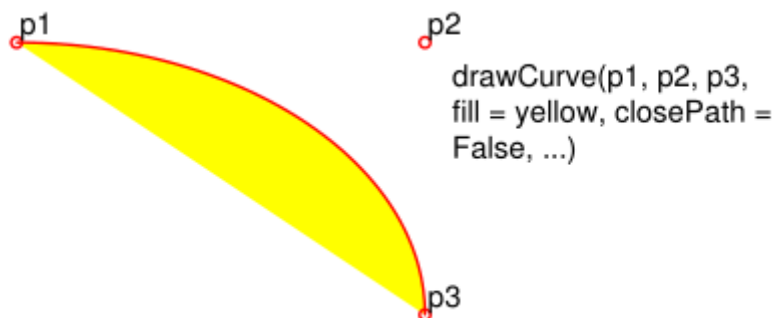
This method may be used to easily draw quadrants of circles or ellipses.

Parameters

- *p1* (*Point*) – is the starting point of the curve.
- *p2* (*Point*) – control point.
- *p3* (*Point*) – is the end point of the curve.

For a description of the other parameters see *Common Parameters*.

`drawCurve()` example:



`drawRect(rect, color = (0, 0, 0), fill = None, width = 1, dashes = None, roundCap = True, overlay = True)`

PDF only: Draw a rectangle. Apart from `rect`, parameters have the same meaning as in method `drawPolyline()`.

Parameters `rect` (*Rect*) – where to put the rectangle on the page. `rect` must be finite and not empty.

For a description of the other parameters see *Common Parameters*.

An efficient way to background-color a PDF page with the old Python paper color is `page.drawRect(page.rect, color = py_color, fill = py_color, overlay = False)`, where `py_color = getColor("py_color")`.

`insertImage(rect, filename = None, pixmap = None, overlay = True)`

PDF only: Fill the given rectangle with an image. Width and height need not have the same proportions as the image: it will be adjusted to fit. The image is either taken from a pixmap or from a file - **exactly one** of these parameters **must** be specified.

Parameters

- `rect` (*Rect*) – where to put the image on the page. `rect` must be finite, not empty and be completely contained in the page's rectangle.
- `filename` (*str*) – name of an image file (all MuPDF supported formats - see *Pixmap* chapter).
- `pixmap` (*Pixmap*) – pixmap containing the image. When inserting the same image multiple times, this should be the preferred option, because the overhead of opening the image and decompressing its content will occur every time with the filename option.

For a description of the other parameters see *Common Parameters*.

Returns zero

This example puts the same image on every page of a document:

```
>>> doc = fitz.open(...)
>>> rect = fitz.Rect(0, 0, 50, 50) # put thumbnail in upper left corner
>>> pix = fitz.Pixmap("some.jpg")  # an image file
>>> for page in doc:
>>>     page.insertImage(rect, pixmap = pix)
>>> doc.save(...)
```

Notes:

1. If that same image had already been present in the PDF, then only a reference will be inserted. This of course considerably saves disk space and processing time. But in order to detect this fact, existing PDF images need to be compared with the new one. This is achieved by storing an md5 code for each image and only compare the new image's code against these. Generating this md5 table, however, is done only when triggered by the first image insertion - which therefore may have an extended response time.

2. You can use this method to provide a background image for the page, like a copyright, a watermark or a background color. Or you can combine it with `searchFor()` to achieve a textmarker effect.
3. The image may be inserted uncompressed, e.g. if a `Pixmap` is used or if the image has an alpha channel. Therefore, consider using `deflate = True` when saving the file.

`getText(output = 'text')`

Retrieves the text of a page. Depending on the output parameter, the results of the [TextPage](#) extract methods are returned.

If 'text' is specified, plain text is returned in the order as specified during PDF creation (which is not necessarily the normal reading order). As this may not always look as expected, consider using the example program `PDF2TextJS.py`. It is based on `output = 'json'` (= `TextPage.extractJSON()`) and re-arranges text according to the Western reading layout convention “from top-left to bottom-right”.

Parameters `output (str)` – A string indicating the requested text format, one of `text` (default), `html`, `json`, or `xml`.

Return type `string`

Returns The page's text as one string.

`getFontList()`

PDF only: Return a list of fonts referenced by the page. Same as [Document.getPageFontList\(\)](#).

`getImageList()`

PDF only: Return a list of images referenced by the page. Same as [Document.getPageImageList\(\)](#).

`getPixmap(matrix = fitz.Identity, colorspace = "RGB", clip = None, alpha = True)`

Creates a Pixmap from the page.

Parameters

- `matrix (Matrix)` – A [Matrix](#) object. Default is the [Identity](#) matrix.
- `colorspace (string, Colorspace)` – Defines the required colorspace, one of `GRAY`, `RGB` (default) or `CMYK` (case insensitive). Alternatively specify a [Colorspace](#), conveniently one of the predefined ones (`csGRAY`, `csRGB` or `csCMYK`).
- `clip (IRect)` – An [IRect](#) to restrict rendering of the page to the rectangle's area. If not specified, the complete page will be rendered.
- `alpha (bool)` – A bool indicating whether an alpha channel should be included in the pixmap. Choose `False` if you do not absolutely need transparency. This will save a lot of memory (25% in case of RGB), and also processing time in most cases.

Return type [Pixmap](#)

Returns Pixmap of the page.

Note: Using `alpha` does not only incur memory, but also processing time penalties: Pixmapes are normally used further on, to .g. either save them as images, or to display them with a GUI manager. In case of [wxPython](#) we have observed, that for pixmaps of the same page image, it is about two times faster to create a `wx.Bitmap` with method `FromBuffer` of a non-alpha version, compared with `FromBufferRGBA` of an alpha pixmap. In combination, using the alpha-free alternative is at least 20% faster (this is true across all Python, [wxPython](#) and [bitness](#) versions). So, our recommendation is to do `alpha = False`, and then `wx.Bitmap.FromBuffer(pix.w, pix.h, pix.samples)`. Similar considerations apply for other GUI managers like `TK` and `Qt`. Processing time difference (alpha vs. non-alpha) within MuPDF itself is below 5%.

`setRotation(rot)`

PDF only: Sets the rotation of the page.

Parameters *rot* (*int*) – An integer specifying the required rotation in degrees. Should be a (positive or negative) multiple of 90.

Returns zero if successful, -1 if not a PDF.

`searchFor(text, hit_max = 16)`

Searches for *text* on a page. Identical to *TextPage.search()*.

Parameters

- *text* (*str*) – Text to searched for. Upper / lower case is ignored.
- *hit_max* (*int*) – Maximum number of occurrences accepted.

Return type list

Returns A list of *Rect* rectangles each of which surrounds one occurrence of *text*.

`rotation`

PDF only: contains the rotation of the page in degrees and -1 for other document types.

Type int

`firstLink`

Contains the first *Link* of a page (or None).

Type *Link*

`firstAnnot`

Contains the first *Annot* of a page (or None).

Type *Annot*

`number`

The page number.

Type int

`parent`

The owning document object.

Type *Document*

4.10.1 Common Parameters

`fontname` (*str*)

In general, there are three options:

1. Use one of the standard *PDF Base 14 Fonts*. In this case, *fontfile* must not be specified and "Helvetica" is used if this parameter is omitted, too.
2. Choose a font already in use by the page. Then specify its name **reference** prefixed with a slash /, see below.
3. Specify a *fontfile* present on your system. In this case choose an arbitrary unique new name for this parameter (without prefix).

If inserted text should re-use one of the page's fonts, use its reference name in *getFontList()* like so:

Suppose the font list has the entry [1024, 0, 'Type1', 'CJXQIC+NimbusMonL-Bold', 'R366'], then specify *fontname* = "/R366", *fontfile* = None in order to use font CJXQIC+NimbusMonL-Bold.

Note: We currently only support single byte characters and horizontal, left-to-right text orientation with our text methods (the `rotate` parameter is not influenced by this). Keep this in mind if you re-use an existing font, or use an external fontfile (next parameter).

fontfile (*str*)

File path of a font existing on your computer. If you specify `fontfile`, make sure you use a `fontname` **not occurring** in the above list. This new font will be embedded in the PDF upon `doc.save()`. Similar to new images, a font will be embedded only once. A table of `md5` codes for the binary font contents is maintained and checked against by MuPDF.

fontsize (*float*)

Font size of text. This also determines the line height as `fontsize * 1.2`.

dashes (*str*)

Causes lines to be dashed. A continuous line with no dashes is drawn with `[]0`, `[]` or `None`. For (the rather complex) details on how to achieve dashing, see [Adobe PDF Reference 1.7](#), page 217. Simple versions may look like `[3 4]`, which means dashes of 3 and gaps of 4 pixels length follow each other. `[3 3]` and `[3]` do the same thing.

color / fill (*list, tuple*)

Line and fill colors are always specified as RGB triples of floats from 0 to 1. To simplify color specification, method `getColor()` in `fitz.utils` may be used. It accepts a string as the name of the color and returns the corresponding triple. The method knows over 540 color names - see section [Color Database](#).

overlay (*bool*)

Causes the item to appear in foreground (default) or background.

roundCap (*bool*)

Cause lines, dashes and edges to be rounded (default). If false, sharp edges and square line and dashes ends will be generated. Rounded lines / dashes will end in a semi-circle with a diameter equal to line width and make longer by the radius of this semi-circle.

closePath (*bool*)

Causes the end point of a drawing to be automatically connected with the starting point (by a straight line).

4.10.2 Description of `getLinks()` Entries

Each entry of the `getLinks()` list is a dictionary with the following keys:

- **kind**: (required) an integer indicating the kind of link. This is one of `LINK_NONE`, `LINK_GOTO`, `LINK_GOTOR`, `LINK_LAUNCH`, or `LINK_URI`. For values and meaning of these names refer to *Enumerations*.
- **from**: (required) a *Rect* describing the “hot spot” location on the page’s visible representation (where the cursor changes to a hand image, usually).
- **page**: a 0-based integer indicating the destination page. Required for `LINK_GOTO` and `LINK_GOTOR`, else ignored.
- **to**: either a `fitz.Point`, specifying the destination location on the provided page, default is `fitz.Point(0, 0)`, or a symbolic (indirect) name. If an indirect name is specified, `page = -1` is required and the name must be defined in the PDF in order for this to work. Required for `LINK_GOTO` and `LINK_GOTOR`, else ignored.
- **file**: a string specifying the destination file. Required for `LINK_GOTOR` and `LINK_LAUNCH`, else ignored.
- **uri**: a string specifying the destination internet resource. Required for `LINK_URI`, else ignored.
- **xref**: an integer specifying the PDF cross reference entry of the link object. Do not change this entry in any way. Required for link deletion and update, otherwise ignored. For non-PDF documents, this entry contains `-1`. It is also `-1` for **all** entries in the `getLinks()` list, if **any** of the links is not supported by MuPDF - see the note below.

4.10.3 Notes on Supporting Links

MuPDF’s support for links has changed in **v1.10a**. These changes affect link types *LINK_GOTO* and *LINK_GOTOR*.

Reading (pertains to method `getLinks()` and the `firstLink` property chain)

If MuPDF detects a link to another file, it will supply either a `LINK_GOTOR` or a `LINK_LAUNCH` link kind. In case of `LINK_GOTOR` destination details may either be given as page number (eventually including position information), or as an indirect destination.

If an indirect destination is given, then this is indicated by `page = -1`, and `link.dest.dest` will contain this name. The dictionaries in the `getLinks()` list will contain this information as the `to` value.

Internal links are always of kind `LINK_GOTO`. If an internal link specifies an indirect destination, it **will always be resolved** and the resulting direct destination will be returned. Names are **never returned for internal links**, and undefined destinations will cause the link to be ignored.

Writing

PyMuPDF writes (updates, inserts) links by constructing and writing the appropriate PDF object **source**. This makes it possible to specify indirect destinations for `LINK_GOTOR` and `LINK_GOTO` link kinds (pre PDF 1.2 file formats are **not supported**).

Caution: If a `LINK_GOTO` indirect destination specifies an undefined name, this link can later on not be found / read again with MuPDF / PyMuPDF. Other readers however **will** detect it, but flag it as erroneous.

Indirect `LINK_GOTOR` destinations can in general of course not be checked for validity and are therefore **always accepted**.

4.10.4 Homologous Methods of Document and Page

This is an overview of homologous methods on the *Document* and on the *Page* level.

| Document Level | Page Level |
|--|--------------------------------------|
| <code>doc.getPageFontlist(pno)</code> | <code>doc[pno].getFontlist()</code> |
| <code>doc.getPageImageList(pno)</code> | <code>doc[pno].getImageList()</code> |
| <code>doc.getPagePixmap(pno, ...)</code> | <code>doc[pno].getPixmap(...)</code> |
| <code>doc.getPageText(pno, ...)</code> | <code>doc[pno].getText(...)</code> |
| <code>doc.searchPageFor(pno, ...)</code> | <code>doc[pno].searchFor(...)</code> |
| <code>doc._getPageXref(pno)</code> | <code>doc[pno]._getXref()</code> |

The list assumes a document object `doc`. The page number `pno` is 0-based and can be any positive or negative number $< \text{len}(\text{doc})$.

4.11 Pixmap

Pixmaps (“pixel maps”) are objects at the heart of MuPDF’s rendering capabilities. They represent plane rectangular sets of pixels. Each pixel is described by a number of bytes (“components”) plus an (optional since v1.10.0) alpha byte.

In PyMuPDF, there exist several ways to create a pixmap. Except one, all of them are available as overloaded constructors. A pixmap can be created ...

1. from a document page (via methods *Page.getPagePixmap()* or *Document.getPagePixmap()*)
2. empty based on *Colorspace* and *IRect* information
3. from an image file
4. from an in-memory image (bytearray)
5. from a memory area of plain pixels
6. from an image inside a PDF document
7. as a copy of another pixmap

Note: A number of image formats is supported as input using the **file** or **in-memory constructors**. For a list see section below.

Have a look at the **example** section to see some pixmap usage “at work”.

| Method / Attribute | Short Description |
|-----------------------------|--|
| <i>Pixmap.clearWith()</i> | clears (parts of) a pixmap |
| <i>Pixmap.copyPixmap()</i> | copy parts of another pixmap |
| <i>Pixmap.gammaWith()</i> | applies a gamma factor to the pixmap |
| <i>Pixmap.getPNGData()</i> | returns a PNG as a memory area |
| <i>Pixmap.invertIRect()</i> | invert the pixels of a given area |
| <i>Pixmap.tintWith()</i> | tints a pixmap with a color |
| <i>Pixmap.writeImage()</i> | saves a pixmap in a variety of image formats |
| <i>Pixmap.writePNG()</i> | saves a pixmap as a PNG file |
| <i>Pixmap.alpha</i> | indicates whether transparency is included |
| <i>Pixmap.colorspace</i> | contains the <i>Colorspace</i> |
| <i>Pixmap.height</i> | height of the region in pixels |
| <i>Pixmap.interpolate</i> | interpolation method indicator |
| <i>Pixmap.irect</i> | is the <i>IRect</i> of the pixmap |
| <i>Pixmap.n</i> | number of bytes per pixel including alpha byte |
| <i>Pixmap.samples</i> | the components data for all pixels |
| <i>Pixmap.size</i> | contains the pixmap's total length |
| <i>Pixmap.stride</i> | number of bytes of one image row |
| <i>Pixmap.width</i> | width of the region in pixels |
| <i>Pixmap.x</i> | X-coordinate of top-left corner of pixmap |
| <i>Pixmap.xres</i> | resolution in X-direction |
| <i>Pixmap.y</i> | Y-coordinate of top-left corner of pixmap |
| <i>Pixmap.yres</i> | resolution in Y-direction |

Class API

class Pixmap

`__init__(self, colorspace, irect, alpha)`

This constructor creates an empty pixmap of a size and an origin specified by the irect object. So, for a `fitz.IRect(x0, y0, x1, y1)`, `fitz.Point(x0, y0)` designates the top left corner of the pixmap. Note that the image area is **not initialized** and will contain crap data.

Parameters

- `colorspace` (*Colorspace*) – The colorspace of the pixmap.
- `irect` (*IRect*) – Specifies the pixmap's area and its location.
- `alpha` (*bool*) – Specifies whether transparency bytes should be included. Default is `False`.

`__init__(self, doc, xref)`

This constructor creates a pixmap with origin (0, 0) from an image contained in PDF document `doc` identified by its XREF number.

Parameters

- `doc` (*Document*) – an opened **PDF** document.
- `xref` (*int*) – the XREF number of the image.

`__init__(self, colorspace, sourcepix[, alpha])`

This constructor creates a new pixmap as a copy of another one, `sourcepix`. If the two colorspace differ, a conversion will take place. Any combination of supported colorspace is possible. If `alpha == 1 / True` (default) the result will have the same alpha as the source, otherwise the target has no alpha channel.

Parameters

- `colorspace` (*Colorspace*) – The colorspace of the pixmap.

- `sourcepixmap (Pixmap)` – the source pixmap.
- `alpha (bool)` – whether to also copy the source’s alpha channel.

`__init__(self, filename)`

This constructor creates a pixmap from the image contained in file `filename`. The image type and all other properties are determined automatically.

Parameters `filename (str)` – Path / name of the file. The origin of the resulting pixmap is (0, 0).

`__init__(self, img)`

This constructor creates a non-empty pixmap from `img`, which is assumed to contain a supported image as a bytearray. The image type and all other properties are determined automatically.

Parameters `img (bytearray)` – Data containing a complete, valid image in one of the supported formats. E.g. this may have been obtained from a statement like `img = bytearray(open('somepic.png', 'rb').read())`. The origin of the resulting pixmap is (0,0).

`__init__(self, colorspace, width, height, samples, alpha)`

This constructor creates a non-empty pixmap from `samples`, which is assumed to contain an image in “plain pixel” format. This means that each pixel is represented by `n` bytes (as controlled by the `colorspace` and `alpha` parameters). The origin of the resulting pixmap is (0,0). This method is useful when raw image data are provided by some other program - see examples below.

Parameters

- `colorspace (Colorspace)` – Colorspace of the image. Together with `alpha` this parameter controls the interpretation of the `samples` area: for `CS_GRAY`, `CS_RGB` and `CS_CMYK`, 1 + `alpha`, 3 + `alpha` or 4 + `alpha` bytes in samples will be assumed to define one pixel, respectively. Calling this number `n`, the following must evaluate to `True`: `n * width * height == len(samples)`.
- `width (int)` – Width of the image in pixels
- `height (int)` – Height of the image in pixels
- `samples (bytearray or bytes or str)` – bytearray, bytes or string (Python 2 only) containing consecutive bytes describing all pixels of the image.
- `alpha (bool)` – a transparency channel is included in samples.

`clearWith(value[, irect])`

Clears an area specified by the `IRect` `irect` within a pixmap. To clear the whole pixmap omit `irect`.

Parameters

- `value (int)` – Values from 0 to 255 are valid. Each color byte of each pixel will be set to this value, while alpha will always be set to 255 (non-transparent) if present. Default is 0 (black).
- `irect (IRect)` – An `IRect` object specifying the area to be cleared.

`tintWith(red, green, blue)`

Colorizes (tints) a pixmap with a color provided as a value triple (red, green, blue). Use this method only for `CS_GRAY` or `CS_RGB` colorspaces. A `TypeError` exception will otherwise be raised.

If the colorspace is `CS_GRAY`, $(red + green + blue)/3$ will be taken as the tinting value.

Parameters

- `red (int)` – The red component. Values from 0 to 255 are valid.

- **green** (*int*) – The **green** component. Values from 0 to 255 are valid.
- **blue** (*int*) – The **blue** component. Values from 0 to 255 are valid.

`gammaWith(gamma)`

Applies a gamma factor to a pixmap, i.e. lightens or darkens it.

Parameters `gamma` (*float*) – `gamma = 1.0` does nothing, `gamma < 1.0` lightens, `gamma > 1.0` darkens the image.

`invertIRect(irect)`

Invert the color of all pixels in an area specified by *IRect* `irect`. To invert everything, use `getIRect()` or omit this parameter.

Parameters `irect` (*IRect*) – The area to be inverted.

`copyPixmap(source, irect)`

Copies the *IRect* part of the **source** pixmap into the corresponding area of this one. The two pixmaps may have different dimensions and different colorspace (provided each is either *CS_GRAY* or *CS_RGB*), but currently **must** have the same alpha property. The copy mechanism automatically adjusts discrepancies between source and target pixmap like so:

If copying from *CS_GRAY* to *CS_RGB*, the source gray-shade value will be put into each of the three rgb component bytes. If the other way round, $(r + g + b) / 3$ will be taken as the gray-shade value of the target.

Between the specified `irect` and the target pixmap's *IRect*, an “intersection” rectangle is calculated at first. Then the corresponding data of this intersection are being copied. If the intersection is empty, nothing will happen.

If you want your **source** pixmap image to land at a specific position of the target, set its **x** and **y** attributes to the top left point of the desired rectangle before copying. See the example below for how this works.

Parameters

- **source** (*Pixmap*) – The pixmap from where to copy.
- **irect** (*IRect*) – An *IRect* object specifying the area to be copied.

`writePNG(filename)`

Saves a pixmap as a PNG file. Please note that only grayscale and RGB colorspace can be saved in PNG format (this is a MuPDF restriction). CMYK colorspace must either be saved as ***.pam** files or be converted. Since MuPDF v1.10a the **savealpha** option is no longer supported and will be ignored with a warning.

Parameters `filename` (*str*) – The filename to save as (the extension **png** must be specified).

`getPNGData()`

Returns the pixmap as an image area (bytearray) in PNG format. Please note that only grayscale and RGB colorspace can be produced in PNG format (this is a MuPDF restriction). CMYK colorspace must be converted first. Since MuPDF v1.10a the **savealpha** option is no longer supported and will be ignored with a warning.

Return type bytearray

`writeImage(filename, output="png")`

Saves a pixmap as an image file. This method is an extension to `writePNG()`. Depending on the output chosen, some or all colorspace are supported and different file extensions can be chosen. Please see the table below. Since MuPDF v1.10a the **savealpha** option is no longer supported and will be ignored with a warning.

Parameters

- **filename** (*str*) – The filename to save to. Depending on the chosen output format, possible file extensions are **.pam**, **.pbm**, **.pgm**, **ppm**, **.pnm**, **.png** and **.tga**.

- **output** (*str*) – The requested image format. The default is `png` for which this function is equivalent to `writePNG()`. Other possible values are `pam`, `pnm` and `tga`.

alpha

Indicates whether this pixmap contains transparency information

Type `bool`

colorspace

The colorspace of the pixmap.

Type `str`

stride

Contains the length of one row of image data in **samples**. This is primarily used for calculation purposes. The following expressions are `True`: `len(samples) == height * stride`, `width * n == stride`, `Colorspace.nbytes + alpha == n`.

Type `int`

irect

Contains the *IRect* of the pixmap.

Type *IRect*

samples

The color and transparency values for all pixels. **samples** is a memory area of size `width * height * n` bytes. Each `n` bytes define one pixel. Each successive `n` bytes yield another pixel in scanline order. Subsequent scanlines follow each other with no padding. E.g. for an RGBA colorspace this means, **samples** is a bytearray like `..., R, G, B, A, ...`, and the four byte values `R, G, B, A` define one pixel.

This area can also be used by other graphics libraries like PIL (Python Imaging Library) to do additional processing like saving the pixmap in other image formats. See example 3.

Type `bytes`

Note: We have changed the type of **samples** from `bytearray` to `bytes`. Some GUIs (Tk) require a read-only type here. We hope this does not interfere with existing code!

size

Contains the total length of the pixmap. This will generally equal `len(pix.samples) + 60`. The following will evaluate to `True`: `len(pixmap) == pixmap.size`.

Type `int`

width

The width of the region in pixels. For compatibility reasons, `w` is also supported.

Type `int`

height

The height of the region in pixels. For compatibility reasons, `h` is also supported.

Type `int`

x

X-coordinate of top-left corner

Type `int`

y

Y-coordinate of top-left corner

Type `int`

n
Number of components per pixel. This number depends on colorspace and alpha (see remark above). `Pixmap.n - Pixmap.alpha == Pixmap.colorspace.n` is always `True`.
Type `int`

xres
Horizontal resolution in dpi (dots per inch).
Type `int`

yres
Vertical resolution in dpi.
Type `int`

interpolate
An information-only boolean flag set to `True` if the image will be drawn using “linear interpolation”. If `False` “nearest neighbour sampling” will be used.
Type `bool`

4.11.1 Supported Pixmap Construction Image Types

The following file types are supported as input to construct pixmaps: **BMP**, **JPEG**, **GIF**, **SVG**, **TIFF**, **JXR**, and **PNG**.

4.11.2 Details on Saving Images with `writeImage()`

The following table shows possible combinations of file extensions, output formats and colorspace of method `writeImage()`:

| output = | CS_GRAY | CS_RGB | CS_CMYK |
|-----------------|----------------|---------------|----------------|
| "pam" | .pam | .pam | .pam |
| "pnm" | .pnm, .pgm | .pnm, .ppm | invalid |
| "png" | .png | .png | invalid |
| "tga" | .tga | .tga | invalid |

Note: Not all image file types are available, or at least common on all platforms, e.g. PAM is mostly unknown on Windows. Especially pertaining to CMYK colorspace, you can always convert a CMYK pixmap to an RGB-pixmap with `rgb_pix = fitz.Pixmap(fitz.csRGB, cmyk_pix)` and then save that as a PNG.

4.11.3 Pixmap Example Code Snippets

Example 1

This shows how pixmaps can be used for purely graphical, non-PDF purposes. The script reads a PNG picture and creates a new PNG file which consist of 3 * 4 tiles of the original one:

```

import fitz
# create a pixmap of a picture
pix0 = fitz.Pixmap("editra.png")

# set target colorspace and pixmap dimensions and create it
tar_width = pix0.width * 3          # 3 tiles per row
tar_height = pix0.height * 4        # 4 tiles per column
tar_irect = fitz.IRect(0, 0, tar_width, tar_height)
# create empty target pixmap
tar_pix = fitz.Pixmap(fitz.csRGB, tar_irect, pix0.alpha)
# clear target with a very lively stone-gray (thanks and R.I.P., Lorient)
tar_pix.clearWith(90)

# now fill target with 3 * 4 tiles of input picture
for i in range(4):
    pix0.y = i * pix0.height        # modify input's y coord
    for j in range(3):
        pix0.x = j * pix0.width      # modify input's x coord
        tar_pix.copyPixmap(pix0, pix0.irect) # copy input to new loc
        # save all intermediate images to show what is happening
        fn = "target-%s-%s.png" % (str(i), str(j))
        tar_pix.writePNG(fn)

```

This is the input picture `editra.png` (taken from the wxPython directory `/tools/Edittra/pixmaps`):



Here is the output, showing some intermediate picture and the final result:



Example 2

This shows how to create a PNG file from a numpy array (several times faster than most other methods):

```
import numpy as np
import fitz
#=====
# create a fun-colored width * height PNG with fitz and numpy
#=====
height = 150
width = 100
bild=np.ndarray((height, width, 3), dtype=np.uint8)

for i in range(height):
    for j in range(width):
        # one pixel (some fun coloring)
        bild[i, j] = [(i+j)%256, i%256, j%256]

samples = bytearray(bild.tostring()) # get plain pixel data from numpy array
pix=fitz.Pixmap(fitz.csRGB, width, height, samples, alpha=False)
pix.writePNG("test.png")
```

Example 3

This shows how to interface with PIL / Pillow (the Python Imaging Library), thereby extending the reach of image files that can be processed:

```
import fitz
from PIL import Image

pix = fitz.Pixmap(...)
... # any code here
# create and save a PIL image
img = Image.frombytes("RGB", [pix.width, pix.height], str(pix.samples))
img.save(filename, 'jpeg')
```



```
# an example for the opposite direction
# create a pixmap from any PIL-supported image file "some_image.xxx"

img = Image.open("some_image.xxx").convert("RGB")
samples = bytearray(img.tobytes())
pix = fitz.Pixmap(fitz.csRGB, img.size[0], img.size[1], samples, alpha=False)
```

4.12 Point

Point represents a point in the plane, defined by its x and y coordinates.

| Attribute / Method | Short Description |
|----------------------------|-------------------------------------|
| <i>Point.distance_to()</i> | calculate distance to point or rect |
| <i>Point.transform()</i> | transform point with a matrix |
| Point.x | the X-coordinate |
| Point.y | the Y-coordinate |

Class API

```
class Point
```

```
__init__(self)
__init__(self, x, y)
__init__(self, point)
__init__(self, list)
```

Overloaded constructors.

Without parameters, `Point(0, 0)` will be created.

With another `point` specified, a **new copy** will be created. A `list` must be Python sequence object of length 2. For a `list`, it is the user's responsibility to only provide numeric entries - **no error checking is done**, and invalid entries will receive a value of `-1.0`.

Parameters

- `x (float)` – X coordinate of the point
- `y (float)` – Y coordinate of the point

```
distance_to(x[, unit])
```

Calculates the distance to `x`, which may be a *Rect*, *IRect* or *Point*. The distance is given in units of either `px` (pixels, default), `in` (inches), `mm` (millimeters) or `cm` (centimeters).

Note: If `x` is a rectangle, the distance is calculated as if the rectangle were finite.

Parameters

- `x (Rect or IRect or Point)` – the object to which the distance is calculated.
- `unit (str)` – the unit to be measured in. One of `px`, `in`, `cm`, `mm`.

Returns distance to object `x`.

Return type float

`transform(m)`

Applies matrix *m* to the point.

Parameters *m* – The matrix to be applied.

Return type Point

4.12.1 Remark

A point's *p* attributes *x* and *y* can also be accessed as indices, e.g. `p.x == p[0]`, and the `tuple()` and `list()` functions yield sequence objects of its components.

4.12.2 Point Algebra

The following arithmetic operators have been defined for the `Point` class.

Binary Operators

- **Addition:** `p + x` is a new `Point` with added coordinates of *p* and *x* (another `Point`, a sequence or a number). If *x* is a number, it is added to both components of *p*.
- **Subtraction:** analogous to addition.
- **Multiplication:** `p * m` is a new `Point` *p* transformed by matrix *m*. If *m* is a number, it will multiply the coordinates.
- **Comparison:** `p1 == p2` is `True` if the coordinates are equal (not only if they are the same object!).

Unary Operators

- `-p` is a copy of *p* with negated coordinates.
- `+p` is a copy of *p*.
- `abs(p)` means the Euclidean norm of *p*, i.e. its length as a vector.
- `bool(p)` is `False` for `Point(0, 0)` and `True` otherwise.

4.12.3 Examples

This should illustrate some basic uses:

```
>>> fitz.Point(1, 2) * fitz.Matrix(90)
fitz.Point(-2.0, 1.0)
>>>
>>> fitz.Point(1, 2) * 3
fitz.Point(3.0, 6.0)
>>>
>>> fitz.Point(1, 2) + 3
fitz.Point(4.0, 5.0)
>>>
>>> fitz.Point(25, 30) + fitz.Point(1, 2)
fitz.Point(26.0, 32.0)
>>> fitz.Point(25, 30) + (1, 2)
fitz.Point(26.0, 32.0)
>>>
```

```
>>> fitz.Point([1, 2])
fitz.Point(1.0, 2.0)
>>>
>>> -fitz.Point(1, 2)
fitz.Point(-1.0, -2.0)
>>>
>>> abs(fitz.Point(25, 30))
39.05124837953327
```

4.13 Rect

Rect represents a rectangle defined by four floating point numbers x_0 , y_0 , x_1 , y_1 . They are viewed as being coordinates of two diagonally opposite points. The first two numbers are regarded as the “top left” corner P_{x_0,y_0} and P_{x_1,y_1} as the “bottom right” one. However, these two properties need not coincide with their ostensive meanings - read on.

The following remarks are also valid for *IRect* objects:

- Rectangle borders are always parallel to the respective X- and Y-axes.
- The constructing points can be anywhere in the plane - they need not even be different, and e.g. “top left” need not be the geometrical “north-western” point.
- For any given quadruple of numbers, the geometrically “same” rectangle can be defined in (up to) four different ways: $\text{Rect}(P_{x_0,y_0}, P_{x_1,y_1})$, $\text{Rect}(P_{x_1,y_1}, P_{x_0,y_0})$, $\text{Rect}(P_{x_0,y_1}, P_{x_1,y_0})$, and $\text{Rect}(P_{x_1,y_0}, P_{x_0,y_1})$.

Hence some useful classification:

- A rectangle is called **finite** if $x_0 \leq x_1$ and $y_0 \leq y_1$ (i.e. the bottom right point is “south-eastern” to the top left one), otherwise **infinite**. Of the four alternatives above, only one is finite (disregarding degenerate cases).
- A rectangle is called **empty** if $x_0 = x_1$ or $y_0 = y_1$, i.e. if its area is zero.

Note: As paradox as it may sound: a rectangle can be both, infinite **and** empty ...

| Methods / Attributes | Short Description |
|----------------------------|---|
| <i>Rect.round()</i> | create smallest <i>IRect</i> containing rectangle |
| <i>Rect.transform()</i> | transform rectangle with a matrix |
| <i>Rect.intersect()</i> | common part with another rectangle |
| <i>Rect.includePoint()</i> | enlarge rectangle to also contain a point |
| <i>Rect.includeRect()</i> | enlarge rectangle to also contain another one |
| <i>Rect.getRectArea()</i> | calculate rectangle area |
| <i>Rect.getArea()</i> | calculate rectangle area |
| <i>Rect.contains()</i> | checks containment of another object |
| <i>Rect.intersects()</i> | checks for non-empty intersections |
| <i>Rect.normalize()</i> | makes a rectangle finite |
| <i>Rect.height</i> | rectangle height |
| <i>Rect.irect</i> | equals result of method <i>round()</i> |
| <i>Rect.width</i> | rectangle width |
| <i>IRect.top_left</i> | top left point |
| <i>IRect.top_right</i> | top_right point |
| <i>IRect.bottom_left</i> | bottom left point |
| <i>IRect.bottom_right</i> | bottom right point |
| <i>Rect.x0</i> | top left corner's X-coordinate |
| <i>Rect.y0</i> | top left corner's Y-coordinate |
| <i>Rect.x1</i> | bottom right corner's X-coordinate |
| <i>Rect.y1</i> | bottom right corner's Y-coordinate |
| <i>Rect.isInfinite</i> | True if rectangle is infinite |
| <i>Rect.isEmpty</i> | True if rectangle is empty |

Class API

class Rect

```
__init__(self)
__init__(self, x0, y0, x1, y1)
__init__(self, top_left, bottom_right)
__init__(self, top_left, x1, y1)
__init__(self, x0, y0, bottom_right)
__init__(self, rect)
__init__(self, list)
```

Overloaded constructors: *top_left*, *bottom_right* stand for *Point* objects, *list* is a Python sequence type with length 4, *rect* means another *Rect*, while the other parameters mean float coordinates. If *list* is specified, it is the user's responsibility to only provide numeric entries - **no error checking is done**, and invalid entries will receive a value of `-1.0`.

If *rect* is specified, the constructor creates a **new copy** of *rect*.

Without any parameters, the rectangle *Rect*(0.0, 0.0, 0.0, 0.0) is created.

round()

Creates the smallest containing *IRect*. This is **not** the same as simply rounding each of the rectangle's coordinates! Look at the example below.

Note: If the rectangle is infinite, the “normalized” (finite) version of it will be rounded. So the result of this method is always a finite *IRect*.

Return type *IRect*

`transform(m)`

Transforms the rectangle with a matrix and **replaces the original**. If the rectangle is empty or infinite, this is a no-operation.

Parameters *m* (*Matrix*) – The matrix for the transformation.

Return type *Rect*

Returns the smallest rectangle that contains the transformed original. This implies that the transformed rectangle's corners in general are **not** the transformed original ones!

Note: If the rectangle is infinite, this is a no-operation. If the rectangle is empty, the result in general is not empty.

`intersect(r)`

The intersection (common rectangular area) of the current rectangle and *r* is calculated and **replaces the current** rectangle. If either rectangle is empty, the result is also empty. If *r* is infinite, this is a no-operation.

Parameters *r* (*Rect*) – Second rectangle

`includeRect(r)`

The smallest rectangle containing the current one and *r* is calculated and **replaces the current** one. If either rectangle is infinite, the result is also infinite. If one is empty, the other one will be taken as the result.

Parameters *r* (*Rect*) – Second rectangle

`includePoint(p)`

The smallest rectangle containing the current one and point *p* is calculated and **replaces the current** one. **Infinite rectangles remain unchanged**. To create a rectangle containing a series of points, start with (the empty) `fitz.Rect(p1, p1)` and successively perform `includePoint` operations for the other points.

Parameters *p* (*Point*) – Point to include.

`getRectArea([unit])`

or

`getArea([unit])`

Calculates the area of the rectangle and, with no parameter, equals `abs(rect)`. Like an empty rectangle, the area of an infinite rectangle is also zero. So, at least one of `fitz.Rect(p1, p2)` and `fitz.Rect(p2, p1)` has a zero area.

Parameters *unit* (*str*) – Specify required unit: respective squares of *px* (pixels, default), *in* (inches), *cm* (centimeters), or *mm* (millimeters).

Return type *float*

`contains(x)`

Checks whether *x* is contained in the rectangle. It may be an *IRect*, *Rect*, *Point* or number. If *x* is an empty rectangle, this is always `True`. Conversely, if the rectangle is empty this is always `False`, if *x* is not an empty rectangle and not a number. If *x* is a number, it will be checked to be one of the four components. *x in rect* and `rect.contains(x)` are equivalent.

Parameters *x* (*IRect* or *Rect* or *Point* or *float*) – the object to check.

Return type *bool*

`intersects(r)`

Checks whether the rectangle and *r* (a *Rect* or *IRect*) have a non-empty rectangle in common. This will always be `False` if either is infinite or empty.

Parameters *r* (*IRect* or *Rect*) – the rectangle to check.

Return type bool

`normalize()`

Replace the rectangle with its finite version. This is done by shuffling the rectangle corners. After completion of this method, the bottom right corner will indeed be south-eastern to the top left one.

`top_left`

Equals `Point(x0, y0)`.

Type *Point*

`top_right`

Equals `Point(x1, y0)`.

Type *Point*

`bottom_left`

Equals `Point(x0, y1)`.

Type *Point*

`bottom_right`

Equals `Point(x1, y1)`.

Type *Point*

`width`

Contains the width of the rectangle. Equals `x1 - x0`.

Return type float

`height`

Contains the height of the rectangle. Equals `y1 - y0`.

Return type float

`x0`

X-coordinate of the left corners.

Type float

`y0`

Y-coordinate of the top corners.

Type float

`x1`

X-coordinate of the right corners.

Type float

`y1`

Y-coordinate of the bottom corners.

Type float

`isInfinite`

True if rectangle is infinite, **False** otherwise.

Type bool

`isEmpty`

True if rectangle is empty, **False** otherwise.

Type bool

4.13.1 Remark

A rectangle's coordinates can also be accessed via index, e.g. `r.x0 == r[0]`, and the `tuple()` and `list()` functions yield sequence objects of its components.

4.13.2 Rect Algebra

The following arithmetic operators have been defined for `Rect` objects (denoted as `r` in the following). Note that in most binary operations, the second operand may also be of type *IRect*, *Point* or numbers.

Binary operators

- **Addition:** `r + x` where `r` is a `Rect` and `x` can be a `Rect`, `IRect`, list / tuple or a number. The result is a new `Rect` with added components of the operands. If `x` is a number, it is added to all components of `r`.
- **Subtraction:** analogous to addition.
- **Inclusion “|”:** `r | x` is the new `Rect` that also includes `x`, which can be an `IRect`, `Rect` or `Point`.
- **Intersection “&”:** `r & x` is a new `Rect` containing the area common to `r` and `x` which can be an `IRect` or `Rect`.
- **Multiplication:** `r * m` is the new `Rect` resulting from `r.transform(m)` for a **matrix** `m` or from multiplication with **number** `m` (coordinate-wise).
- **Containment Test:** `if x in r:` tests whether `r` contains `x`. For a `Rect` or `IRect` this tests whether its area is contained in `r`. If `x` is a number, it tests whether `x` is one of the 4 coordinates.
- **Comparison:** `r1 == r2` is `True` if their tuples of coordinates are equal (not only if they are the same object!). This tests floats for equality, so be wary of artifact differences. However, `rect == irect` is always `False` because of the different object types.

Unary Operators

- `-r` is a new copy of `r` with negated components.
- `+r` is a new copy of `r`.
- `bool(r)` is `False` for `Rect(0, 0, 0, 0)` and `True` otherwise.
- `abs(r)` is equal to `r.getArea()`.

4.13.3 Examples

Example 1:

```
>>> p1 = fitz.Point(10, 10)
>>> p2 = fitz.Point(300, 450)
>>>
>>> fitz.Rect(p1, p2)
fitz.Rect(10.0, 10.0, 300.0, 450.0)
>>>
>>> fitz.Rect(10, 10, 300, 450)
fitz.Rect(10.0, 10.0, 300.0, 450.0)
>>>
>>> fitz.Rect(10, 10, p2)
fitz.Rect(10.0, 10.0, 300.0, 450.0)
>>>
>>> fitz.Rect(p1, 300, 450)
fitz.Rect(10.0, 10.0, 300.0, 450.0)
```

Example 2:

```
>>> r = fitz.Rect(0.5, -0.01, 123.88, 455.123456)
>>>
>>> r
fitz.Rect(0.5, -0.009999999776482582, 123.87999725341797, 455.1234436035156)
>>>
>>> r.round()      # = r.irect
fitz.IRect(0, -1, 124, 456)
```

Example 3:

```
>>> m = fitz.Matrix(45)
>>> r = fitz.Rect(10, 10, 410, 610)
>>> r * m
fitz.Rect(-424.2640686035156, 14.142135620117188, 282.84271240234375, 721.2489013671875)
>>>
>>> r | fitz.Point(5, 5)
fitz.Rect(5.0, 5.0, 410.0, 610.0)
>>>
>>> r + 5
fitz.Rect(15.0, 15.0, 415.0, 615.0)
>>>
>>> r & fitz.Rect(0, 0, 15, 15)
fitz.Rect(10.0, 10.0, 15.0, 15.0)
```

Example 4:

```
>>> r = fitz.Rect(...)      # any rectangle
>>> ir = r.irect             # its IRect version
>>> # even though you get ...
>>> ir in r
True
>>> # ... and ...
>>> r in ir
True
>>> # ... r and ir are still different types!
>>> r == ir
False
>>> # corners are always part of non-empty rectangles
>>> r.bottom_left in r
True
>>>
>>> # numbers are checked against coordinates
>>> r.x0 in r
True
```

Example 5:

Create a copy that is **guaranteed to be finite** in two ways:

```
>>> r = fitz.Rect(...)      # any rectangle
>>>
>>> # alternative 1
>>> s = fitz.Rect(r.top_left, r.top_left)  # just a point
>>> s | r.bottom_right      # s is a finite rectangle!
>>>
>>> # alternative 2
>>> s = (+r).normalize()
```


LOW LEVEL FUNCTIONS AND CLASSES

Contains a number of functions and classes for the experienced user. To be used for special needs or performance requirements.

5.1 Functions

The following are miscellaneous functions to be used by the experienced PDF programmer.

| Function | Short Description |
|-------------------------------------|---|
| <i>Annot._cleanContents()</i> | PDF only: clean the annot's /Contents objects |
| <i>Annot._getXref()</i> | PDF only: return XREF number of annotation |
| <i>Document._delXmlMetadata()</i> | PDF only: remove XML metadata |
| <i>Document._getCharWidths()</i> | PDF only: return a list of glyph widths of a font |
| <i>Document._getNewXref()</i> | PDF only: create and return a new XREF entry |
| <i>Document._getObjectString()</i> | PDF only: return object source code |
| <i>Document._getOLRootNumber()</i> | PDF only: return / create XREF of /Outline |
| <i>Document._getPageObjNumber()</i> | PDF only: return XREF and generation number of a page |
| <i>Document._getPageRectText()</i> | PDF only: return raw string within rectangle |
| <i>Document._getPageXref()</i> | PDF only: same as <i>_getPageObjNumber()</i> |
| <i>Document._getXrefLength()</i> | PDF only: return length of XREF table |
| <i>Document._getXrefStream()</i> | PDF only: return content of a stream |
| <i>Document._getXrefString()</i> | PDF only: return object source code |
| <i>Document._updateObject()</i> | PDF only: insert or update a PDF object |
| <i>Document._updateStream()</i> | PDF only: replace the stream of an object |
| <i>getPDFnow()</i> | return the current timestamp in PDF format |
| <i>getPDFstr()</i> | return PDF-compatible string |
| <i>Page._cleanContents()</i> | PDF only: clean the page's /Contents objects |
| <i>Page._getContents()</i> | PDF only: return a list of content numbers |
| <i>Page._getRectText()</i> | PDF only: return raw string within rectangle |
| <i>Page._getXref()</i> | PDF only: return XREF number of page |
| <i>Page.run()</i> | run a page through a device |
| <i>PaperSize()</i> | return width, height for known paper formats |

PaperSize(s)

Convenience function to return width and height of a known paper format code. These values are given in pixels for the standard resolution 72 pixels = 1 inch.

Currently defined formats include A0 through A10, B0 through B10, C0 through C10, Card-4x6, Card-5x7, Commercial, Executive, Invoice, Ledger, Legal, Legal-13, Letter, Monarch and Tabloid-Extra.

A format name must be supplied as a string (case insensitive), optionally suffixed with "L" (landscape) or "P" (portrait). No suffix defaults to portrait.

Parameters *s* (*str*) – a format name like "A4" or "letter-1".

Return type tuple

Returns (width, height) of the paper format. For an unknown format (-1, -1) is returned. `PaperSize("A4")` returns (595, 842) and `PaperSize("letter-1")` delivers (792, 612).

`getPDFnow()`

Convenience function to return the current local timestamp in PDF compatible format, e.g. `D:20170501121525-04'00'` for local datetime May 1, 2017, 12:15:25 in a timezone 4 hours westward of the UTC meridian.

Return type str

Returns current local PDF timestamp.

`getPDFstr(obj, brackets = True)`

Make a PDF-compatible string: if `obj` contains code points `ord(c) > 255`, then it will be converted to UTF-16BE as a hexadecimal character string like `<feff...>`. Otherwise, if `brackets = True`, it will enclose the argument in `()` replacing any characters with code points `ord(c) > 127` by their octal number `\nnn` prefixed with a backslash. If `brackets = False`, then the string is returned unchanged.

Parameters `obj` (*str or bytes or unicode*) – the object to convert

Return type str

Returns PDF-compatible string enclosed in either `()` or `<>`.

`Document._delXmlMetadata()`

PDF documents only: Delete an object containing XML-based metadata from the PDF. (Py-) MuPDF does not support XML-based metadata. Use this if you want to make sure that the conventional metadata dictionary will be used exclusively. Many third-party PDF programs insert their own metadata in XML format and thus may override what you store in the conventional dictionary. This deletes any such reference, and the corresponding PDF object will be deleted during next garbage collection of the file.

`Document._getPageObjNumber(pno)`

or

`Document._getPageXref(pno)`

PDF documents only: Return the XREF and generation number for a given page.

Parameters `pno` (*int*) – Page number (zero-based).

Return type list

Returns XREF and generation number of page `pno` as a list `[xref, gen]`.

`Page._getXref()`

PDF documents only: Page version for `_getPageObjNumber()` only delivering the XREF (not the generation number).

`Page.run(dev, transform)`

Run a page through a device.

Parameters

- `dev` (*Device*) – Device, obtained from one of the *Device* constructors.
- `transform` (*Matrix*) – Transformation to apply to the page. Set it to *Identity* if no transformation is desired.

`Page._getContents()`

PDF documents only: Return a list of XREF numbers of `/Contents` objects associated with the page.

Return type list

Returns a list of integers, each pointing to an XREF of a `/Contents` object.

Each page has one or more associated contents objects (streams) which contain PDF operator syntax describing what appears where on the page (like text or images, etc. See the *Adobe PDF Reference 1.7*, chapter “Operator Summary”, page 985). This function only enumerates the XREF number(s) of such objects. To get the actual stream source, use function `Document._getXrefStream()` with one of the numbers in this list. Use `Document._updateStream()` to replace the content¹².

`Page._cleanContents()`

PDF documents only: Clean all `/Contents` objects associated with this page (including contents of all annotations). “Cleaning” includes syntactical corrections, standardizations and “pretty printing” of the contents stream. If a page has several contents objects, they will be combined into one. Any discrepancies between `/Contents` and `/Resources` objects are also resolved / corrected. Note that the resulting contents stream will be stored uncompressed (if you do not specify `deflate` on save). See `Page._getContents()` for more details.

Return type int

Returns 0 if successfull (exception is else raised).

`Annot._getXref()`

PDF documents only: Return the xref number of an annotation.

Return type int

Returns XREF number of the annotation.

`Annot._cleanContents()`

PDF documents only: Clean the `/Contents` streams associated with the annotation. This is the same type of action `Page._cleanContents()` performs - just restricted to this annotation.

Return type int

Returns 0 if successful (exception raised otherwise).

`Document._getCharWidths(fontname = None, fontfile = None, xref = 0, limit = 256)`

PDF documents only: Return a list of character (glyph) widths for a font. A font must be specified by exactly one of the parameters `fontname`, `fontfile` or `xref`.

Parameters

- `fontname (str)` – name of a *PDF Base 14 Fonts*. Excludes parameters `fontfile` and `xref`.
- `fontfile (str)` – path / name of a font file available on your system. Excludes parameters `fontname` and `xref`.
- `xref (int)` – cross reference number of a font embedded in the PDF. Excludes parameters `fontname` and `fontfile`. To find a font xref, use e.g. `doc.getPageFontList(pno)` of page number `pno` and take the first entry of one of the returned list entries.
- `limit (int)` – limits the number of returned entries. The default of 256 is sufficient for all fonts that only support characters up to unicode point 255. Specify a number as required.

¹ If a page has multiple contents streams, they are treated as being one logical stream when the document is processed by reader software. A single operator cannot be split between stream boundaries, but a single **instruction** may well be! E.g. for invoking the display of an image, a complete instruction may look like `q a b c d e f cm /imageid Do Q`. In this example, any single of these items (PDF notation: “lexical tokens”) is completely contained in one stream, but `q a b c d e f cm` may be in one and `/imageid Do Q` may be in the next one.

² Note that `/Contents` objects (similar to `/Resources`) may be **shared** among pages. If you change a contents stream, this will affect all pages referencing the same object. To avoid this, use `Page._cleanContents()` **before** making your changes.

Return type list

Returns a list of limit floats, each representing the horizontal width in pixels, that a character needs which has a unicode point equal to an index entry. In order to get the actual width of some character “c”, use `widthlist[ord(c)] * fontsize`. Currently, only horizontal spacing is supported. A zero entry in this list indicates, that the font does not support this unicode point with a glyph. It is up to you to take appropriate action in such cases. Many fonts will have zero entries for indices < 32 (which represents the space character 0x20), others only provide glyphs for the ASCII character set.

A fairly simple function can be used to calculate the pixel width of a string named `text`, like so:

```
def pixlen(text, widthlist, fontsize):
    try:
        return sum([widthlist[ord(c)] for c in text]) * fontsize
    except IndexError:
        m = max([ord(c) for c in text])
        raise ValueError("max. code point found: %i, increase limit" % m)
```

`Document._getPageRectText(pno, rect)`

PDF documents only: Return raw text contained in a rectangle.

Parameters

- `pno` (*int*) – page number.
- `rect` (*Rect*) – rectangle

Return type string

Returns text contained in the rectangle

`Page._getRectText(rect)`

PDF documents only: Return raw plain text contained in a rectangle.

Parameters `rect` (*Rect*) – rectangle

Return type string

Returns text contained in the rectangle

`Document._getObjectString(xref)`

or

`Document._getXrefString(xref)`

PDF documents only: Return the string (“source code”) representing an arbitrary object. For stream objects, only the non-stream part is returned. To get the stream content, use `_getXrefStream()`.

Parameters `xref` (*int*) – XREF number.

Return type string

Returns the string defining the object identified by `xref`.

`Document._getNewXref()`

PDF documents only: Increase the XREF by one entry and return that number. This can then be used to insert a new object.

Return type int

Returns the number of the new XREF entry.

`Document._updateObject(xref, obj_str, page = None)`

PDF documents only: Associate the object identified by string `obj_str` with the XREF number `xref`, which must already exist. If `xref` pointed to an existing object, this will

be replaced with the new object. If a page object is specified, links and other annotations of this page will be reloaded after the object has been updated.

Parameters

- `xref` (*int*) – XREF number.
- `obj_str` (*str*) – a string containing a valid PDF object definition.
- `page` (*Page*) – a page object. If provided, indicates, that annotations of this page should be refreshed (reloaded) to reflect changes incurred with links and / or annotations.

Return type `int`

Returns zero if successful, otherwise an exception will be raised.

`Document._getXrefLength()`

PDF documents only: Return length of XREF table.

Return type `int`

Returns the number of entries in the XREF table.

`Document._getXrefStream(xref)`

PDF documents only: Return decompressed content stream of the object referenced by `xref`. If the object has / is no stream, an exception is raised.

Parameters `xref` (*int*) – XREF number.

Return type `str` or `bytes`

Returns the (decompressed) stream of the object. This is a string in Python 2 and a `bytes` object in Python 3.

`Document._updateStream(xref, stream)`

PDF documents only: Replace the stream of an object identified by `xref`. If the object has no stream, an exception is raised. The function automatically performs a compress operation (“deflate”).

Parameters

- `xref` (*int*) – XREF number.
- `stream` (*bytes* or *bytearray*) – the new content of the stream.

Return type `int`

This method is intended to manipulate streams containing PDF operator syntax (see pp. 985 of the *Adobe PDF Reference 1.7*) as it is the case for e.g. page content streams.

If you update a contents stream, you should use save parameter `clean = True`. This ensures consistency between PDF operator source and the object structure.

Example: Let us assume that you no longer want a certain image appear on a page. This can be achieved by deleting² the respective reference in its contents source(s) - and indeed: the image will be gone after reloading the page. But the page’s `/Resources` object would still³ show the image as being referenced by the page. This save option will clean up any such mismatches.

`Document._getOLRootNumber()`

PDF documents only: Return XREF number of the `/Outlines` root object (this is **not** the first outline entry!). If this object does not exist, a new one will be created.

Return type `int`

Returns XREF number of the `/Outlines` root object.

³ Resources objects are inheritable. This means that many pages can share one. Keeping a page’s `/Resources` object in sync with changes of its `/Contents` therefore may require creating an own `/Resources` object for the page. This can be achieved by either specifying the `clean` option when saving, or by invoking `Page._cleanContents()`.

5.2 Device

The different format handlers (pdf, xps, etc.) interpret pages to a “device”. These devices are the basis for everything that can be done with a page: rendering, text extraction and searching. The device type is determined by the selected construction method.

Class API

class Device

`__init__(self, object, clip)`

Constructor for either a pixel map or a display list device.

Parameters

- `object` (*Pixmap* or *DisplayList*) – one of Pixmap or DisplayList
- `clip` (*IRect*) – An optional *IRect* for Pixmap devices only to restrict rendering to a certain area of the page. If the complete page is required, specify `None`. For display list devices, this parameter must be omitted.

`__init__(self, textsheet, textpage)`

Constructor for a text page device.

Parameters

- `textsheet` (*TextSheet*) – TextSheet object
- `textpage` (*TextPage*) – TextPage object

5.3 DisplayList

DisplayList is a list containing drawing commands (text, images, etc.). The intent is two-fold:

1. as a caching-mechanism to reduce parsing of a page
2. as a data structure in multi-threading setups, where one thread parses the page and another one renders pages.

A **DisplayList** is populated with objects from a page by running *Page.run()* on a *Device*. Replay the list (once or many times) by invoking the display list’s *run()* function.

| Method | Short Description |
|--------------|---|
| <i>run()</i> | (Re)-run a display list through a device. |

Class API

class DisplayList

`__init__(self, mediabox)`

Create a new display list.

When the device is rendering a page, it will populate the display list with drawing commands (text, images, etc.). The display list can later be reused to render a page many times without having to re-interpret the page from the document file.

Parameters `mediabox` (*Rect*) – The page’s rectangle - output of `page.bound()`.

Return type DisplayList

`run(self, dev, ctm, area)`

Parameters

- **dev** (*Device*) – Device
- **ctm** (*Matrix*) – Transformation matrix to apply to display list contents.
- **area** (*Rect*) – Only the part of the contents of the display list visible within this area will be considered when the list is run through the device. This does not apply for tile objects contained in the display list.

5.4 TextPage

TextPage represents the text of a page.

| Method | Short Description |
|-------------------------------|--|
| <i>TextPage.extractText()</i> | Extract the page's plain text |
| <i>TextPage.extractHTML()</i> | Extract the page's text in HTML format |
| <i>TextPage.extractJSON()</i> | Extract the page's text in JSON format |
| <i>TextPage.extractXML()</i> | Extract the page's text in XML format |
| <i>TextPage.search()</i> | Search for a string in the page |

Class API

```
class TextPage
```

extractText()

Extract the text from a **TextPage** object. Returns a string of the page's complete text. No attempt is being made to adhere to a natural reading sequence: the text is returned UTF-8 encoded and in the same sequence as the PDF creator specified it. If this looks awkward for your PDF file, consider using program that re-arranges the text according to a more familiar layout, e.g. `PDF2TextJS.py` in the examples directory.

Return type str

extractHTML()

Extract the text from a **TextPage** object in HTML format. This version contains some more formatting information about how the text is being displayed on the page. See the tutorial chapter for an example.

Return type str

extractJSON()

Extract the text from a **TextPage** object in JSON format. This version contains significantly more formatting information about how the text is being displayed on the page. It is almost as complete as the **extractXML** version, except that positioning information is detailed down to the span level, not to a single character. See the tutorial chapter for an example. To process the returned JSON text use one of the json modules like `json`, `simplejson`, `ujson`, `cjson`, etc. See example program `PDF2TextJS.py` for how to do that.

Return type str

extractXML()

Extract the text from a **TextPage** object in XML format. This contains complete formatting information about every single text character on the page: font, size, line, paragraph, location, etc. This may easily reach several hundred kilobytes of uncompressed data for a text oriented page. See the tutorial chapter for an example.

Return type str

search(string, hit_max = 16)

Search for **string**.

Parameters

- `string (str)` – The string to search for.
- `hit_max (int)` – Maximum number of expected hits (default 16).

Return type list

Returns a list of [Rect](#) objects (without transformation), each surrounding a found string occurrence.

Note: All of the above can be achieved by using the appropriate [Document.getPageText\(\)](#), [Page.getText\(\)](#) and [Page.searchFor\(\)](#) methods.

5.5 TextSheet

`TextSheet` contains a list of distinct text styles used on a page (or a series of pages).

5.6 Working together: Device, DisplayList, TextPage and TextSheet

Here are some instructions on how to use these classes together.

In some situations, performance improvements may be achievable when you fall back to this detail. This is possible when several different things need to be done with the same page - as is demonstrated in the following overview.

5.6.1 Generate Pixmap

The following creates a Pixmap from a document's page (all of this happens behind the curtain when you use method `page.getPixmap()` ...!):

```
mediabox = page.rect                                # save page rectangle for later
dl = fitz.DisplayList(mediabox)                     # requires the mediabox
page.run(fitz.Device(dl), fitz.Identity)             # run page thru a draw Device
rect = (+mediabox).transform(matrix)                 # get the transformed mediabox
pix = fitz.Pixmap(fitz.csRGB, rect.irect)            # (1) allocate an RGB pixmap
pix.clearWith(255)                                  # (2) init it with "white"
# now fill it with the image of the page:
dl.run(fitz.Device(pix, None), fitz.Identity, rect)  # (3) uff!
# but in order to create another pixmap, just re-execute lines (1) thru (3)
```

5.6.2 Perform Text Search

With the existing objects from above, create a new text Device to search for a text string on the page. The `page` object itself is no longer needed (it could have been set to `None`).

For this we need to create `TextPage` and `TextSheet` objects:

```
ts = fitz.TextSheet()
tp = fitz.TextPage(mediabox)                        # TextPage needs original mediabox
dl.run(fitz.Device(ts, tp), fitz.Identity, rect)    # run DisplayList thru text Device
rlist = tp.search("needle")                         # look up "needle"
for r in rlist:                                     # to mark found locations:
    pix.invertIRect(r.irect)                        # invert rectangle colors
```


5.6.3 Extract Text

Again with the existing objects, we can immediately use one or all of the 4 text extraction methods:

```
txt = tp.extractText()      # plain text of the page
json = tp.extractJSON()    # json format
html = tp.extractHTML()    # HTML format
xml = tp.extractXML()      # XML format
```


CONSTANTS AND ENUMERATIONS

Constants and enumerations of MuPDF as implemented by PyMuPDF. Each of the following variables is accessible as `fitz.variable`.

6.1 Constants

Base14_Fonts

Predefined Python list of valid *PDF Base 14 Fonts*.

Return type list

csRGB

Predefined RGB colorspace `fitz.Colorspace(fitz.CS_RGB)`.

Return type *Colorspace*

csGRAY

Predefined GRAY colorspace `fitz.Colorspace(fitz.CS_GRAY)`.

Return type *Colorspace*

csCMYK

Predefined CMYK colorspace `fitz.Colorspace(fitz.CS_CMYK)`.

Return type *Colorspace*

CS_RGB

1 - Type of *Colorspace* is RGBA

Return type int

CS_GRAY

2 - Type of *Colorspace* is GRAY

Return type int

CS_CMYK

3 - Type of *Colorspace* is CMYK

Return type int

VersionBind

'x.xx.x' - version of PyMuPDF (these bindings)

Return type string

VersionFitz

'x.xxx' - version of MuPDF

Return type string

VersionDate

ISO timestamp YYYY-MM-DD HH:MM:SS when these bindings were built.

Return type string

Note: The docstring of `fitz` contains information of the above which can be retrieved like so: `print(fitz.__doc__)`, and should look like: PyMuPDF 1.10.0: Python bindings for the MuPDF 1.10 library, built on 2016-11-30 13:09:13.

6.2 Text Alignment

TEXT_ALIGN_LEFT

0 - align left.

TEXT_ALIGN_CENTER

1 - align center.

TEXT_ALIGN_RIGHT

2 - align right.

TEXT_ALIGN_JUSTIFY

3 - align justify.

6.3 Enumerations

Possible values of `linkDest.kind` (link destination kind). For details consult *Adobe PDF Reference 1.7*, chapter 8.2 on pp. 581.

LINK_NONE

0 - No destination. Indicates a dummy link.

Return type int

LINK_GOTO

1 - Points to a place in this document.

Return type int

LINK_URI

2 - Points to a URI - typically a resource specified with internet syntax.

Return type int

LINK_LAUNCH

3 - Launch (open) another file (of any “executable” type).

Return type int

LINK_GOTOR

5 - Points to a place in another PDF document.

Return type int

6.4 Link Destination Flags

Note: The rightmost byte of this integer is a bit field, so test the truth of these bits with the `&` operator.

LINK_FLAG_L_VALID

1 (bit 0) Top left x value is valid

Return type bool

LINK_FLAG_T_VALID

2 (bit 1) Top left y value is valid

Return type bool

LINK_FLAG_R_VALID

4 (bit 2) Bottom right x value is valid

Return type bool

LINK_FLAG_B_VALID

8 (bit 3) Bottom right y value is valid

Return type bool

LINK_FLAG_FIT_H

16 (bit 4) Horizontal fit

Return type bool

LINK_FLAG_FIT_V

32 (bit 5) Vertical fit

Return type bool

LINK_FLAG_R_IS_ZOOM

64 (bit 6) Bottom right x is a zoom figure

Return type bool

6.5 Annotation Types

Possible values (integer) for PDF annotation types. See chapter 8.4.5, pp. 615 of the Adobe manual for more details.

ANNOT_TEXT

0 - Text annotation

ANNOT_LINK

1 - Link annotation

ANNOT_FREETEXT

2 - Free text annotation

ANNOT_LINE

3 - Line annotation

ANNOT_SQUARE

4 - Square annotation

ANNOT_CIRCLE

5 - Circle annotation

ANNOT_POLYGON

6 - Polygon annotation

ANNOT_POLYLINE

7 - PolyLine annotation

ANNOT_HIGHLIGHT

8 - Highlight annotation

ANNOT_UNDERLINE

9 - Underline annotation

ANNOT_SQUIGGLY
10 - Squiggly-underline annotation

ANNOT_STRIKEOUT
11 - Strikeout annotation

ANNOT_STAMP
12 - Rubber stamp annotation

ANNOT_CARET
13 - Caret annotation

ANNOT_INK
14 - Ink annotation

ANNOT_POPUP
15 - Pop-up annotation

ANNOT_FILEATTACHMENT
16 - File attachment annotation

ANNOT_SOUND
17 - Sound annotation

ANNOT_MOVIE
18 - Movie annotation

ANNOT_WIDGET
19 - Widget annotation

ANNOT_SCREEN
20 - Screen annotation

ANNOT_PRINTERMARK
21 - Printers mark annotation

ANNOT_TRAPNET
22 - Trap network annotation

ANNOT_WATERMARK
23 - Watermark annotation

ANNOT_3D
24 - 3D annotation

6.6 Annotation Flags

Possible mask values for PDF annotation flags.

Note: Annotation flags is a bit field, so test the truth of its bits with the `&` operator. When changing flags for an annotation, use the `|` operator to combine several values. The following descriptions were extracted from the Adobe manual, pages 608 pp.

ANNOT_XF_Invisible
1 - If set, do not display the annotation if it does not belong to one of the standard annotation types and no annotation handler is available. If clear, display such an unknown annotation using an appearance stream specified by its appearance dictionary, if any.

ANNOT_XF_Hidden
2 - If set, do not display or print the annotation or allow it to interact with the user, regardless of its annotation type or whether an annotation handler is available. In cases where screen space is limited, the ability to hide and show annotations selectively can be used in combination with

appearance streams to display auxiliary pop-up information similar in function to online help systems.

ANNOT_XF_Print

4 - If set, print the annotation when the page is printed. If clear, never print the annotation, regardless of whether it is displayed on the screen. This can be useful, for example, for annotations representing interactive pushbuttons, which would serve no meaningful purpose on the printed page.

ANNOT_XF_NoZoom

8 - If set, do not scale the annotation's appearance to match the magnification of the page. The location of the annotation on the page (defined by the upper-left corner of its annotation rectangle) remains fixed, regardless of the page magnification.

ANNOT_XF_NoRotate

16 - If set, do not rotate the annotation's appearance to match the rotation of the page. The upper-left corner of the annotation rectangle remains in a fixed location on the page, regardless of the page rotation.

ANNOT_XF_NoView

32 - If set, do not display the annotation on the screen or allow it to interact with the user. The annotation may be printed (depending on the setting of the Print flag) but should be considered hidden for purposes of on-screen display and user interaction.

ANNOT_XF_ReadOnly

64 - If set, do not allow the annotation to interact with the user. The annotation may be displayed or printed (depending on the settings of the NoView and Print flags) but should not respond to mouse clicks or change its appearance in response to mouse motions.

ANNOT_XF_Locked

128 - If set, do not allow the annotation to be deleted or its properties (including position and size) to be modified by the user. However, this flag does not restrict changes to the annotation's contents, such as the value of a form field.

ANNOT_XF_ToggleNoView

256 - If set, invert the interpretation of the NoView flag for certain events. A typical use is to have an annotation that appears only when a mouse cursor is held over it.

ANNOT_XF_LockedContents

512 - If set, do not allow the contents of the annotation to be modified by the user. This flag does not restrict deletion of the annotation or changes to other annotation properties, such as position and size.

6.7 Annotation Line End Styles

The following descriptions are taken from the Adobe manual TABLE 8.27 on page 630.

ANNOT_LE_None

0 - No line ending.

ANNOT_LE_Square

1 - A square filled with the annotation's interior color, if any.

ANNOT_LE_Circle

2 - A circle filled with the annotation's interior color, if any.

ANNOT_LE_Diamond

3 - A diamond shape filled with the annotation's interior color, if any.

ANNOT_LE_OpenArrow

4 - Two short lines meeting in an acute angle to form an open arrowhead.

`ANNOT_LE_ClosedArrow`

5 - Two short lines meeting in an acute angle as in the `OpenArrow` style (see above) and connected by a third line to form a triangular closed arrowhead filled with the annotation's interior color, if any.

`ANNOT_LE_Butt`

6 - (PDF 1.5) A short line at the endpoint perpendicular to the line itself.

`ANNOT_LE_ROpenArrow`

7 - (PDF 1.5) Two short lines in the reverse direction from `OpenArrow`.

`ANNOT_LE_RClosedArrow`

8 - (PDF 1.5) A triangular closed arrowhead in the reverse direction from `ClosedArrow`.

`ANNOT_LE_Slash`

9 - (PDF 1.6) A short line at the endpoint approximately 30 degrees clockwise from perpendicular to the line itself.

COLOR DATABASE

Since the introduction of methods involving colors (like `Page.drawCircle()`), a requirement may be to have access to predefined colors.

The fabulous GUI package `wxPython` has a database of over 540 predefined RGB colors, which are given more or less memorable names. Among them are not only standard names like “green” or “blue”, but also “turquoise”, “skyblue”, and 100 (not only 50 ...) shades of “gray”, etc.

We have taken the liberty to take a copy of this database (it actually is a list of tuples) modified into PyMuPDF and make its colors available as PDF compatible float triples: for `wxPython`’s (“WHITE”, 255, 255, 255) we return (1, 1, 1), which can be directly used in `color` and `fill` parameters. We also accept any mixed case of “wHiTe” to find a color.

7.1 Function `getColor()`

As the color database may not be needed very often, one additional import statement seems acceptable to get access to it:

```
>>> # "getColor" is the only method you really need
>>> from fitz.utils import getColor
>>> getColor("aliceblue")
(0.9411764705882353, 0.9725490196078431, 1.0)
>>> #
>>> # to get a list of all existing names
>>> from fitz.utils import getColorList
>>> cl = getColorList()
>>> cl
['ALICEBLUE', 'ANTIQUEWHITE', 'ANTIQUEWHITE1', 'ANTIQUEWHITE2', 'ANTIQUEWHITE3',
'ANTIQUEWHITE4', 'AQUAMARINE', 'AQUAMARINE1'] ...
>>> #
>>> # to see the full integer color coding
>>> from fitz.utils import getColorInfoList
>>> il = getColorInfoList()
>>> il
[('ALICEBLUE', 240, 248, 255), ('ANTIQUEWHITE', 250, 235, 215),
('ANTIQUEWHITE1', 255, 239, 219), ('ANTIQUEWHITE2', 238, 223, 204),
('ANTIQUEWHITE3', 205, 192, 176), ('ANTIQUEWHITE4', 139, 131, 120),
('AQUAMARINE', 127, 255, 212), ('AQUAMARINE1', 127, 255, 212)] ...
```

7.2 Printing the Color Database

If you want to actually see how the many available colors look like, use scripts `colordbRGB.py` or `colordbHSV.py` in the examples directory. They create PDFs (already existing in the same directory) with all these colors. Their only difference is sorting order: one takes the RGB values, the other one the Hue-Saturation-Values as sort criteria. This is a screen print of what these files look like.



APPENDIX 1: PERFORMANCE

We have tried to get an impression on PyMuPDF's performance. While we know this is very hard and a fair comparison is almost impossible, we feel that we at least should provide some quantitative information to justify our bold comments on MuPDF's **top performance**.

Following are three sections that deal with different aspects of performance:

- document parsing
- text extraction
- image rendering

In each section, the same fixed set of PDF files is being processed by a set of tools. The set of tools varies - for reasons we will explain in the section.

Here is the list of files we are using. Each file name is accompanied by further information: **size** in bytes, number of **pages**, number of bookmarks (**toc** entries), number of **links**, **text** size as a percentage of file size, **KB** per page, PDF **version** and remarks. **text %** and **KB index** are indicators for whether a file is text or graphics oriented.

| name | size | pages | toc size | links | text % | KB index | version | remarks |
|-----------------|------------|-------|----------|--------|--------|----------|---------|--|
| Adobe.pdf | 32.472.771 | 1.310 | 794 | 32.096 | 8,0% | 24 | PDF 1.6 | linearized, text oriented, many links / bookmarks |
| Evolution.pdf | 13.497.490 | 75 | 15 | 118 | 1,1% | 176 | PDF 1.4 | graphics oriented |
| PyMuPDF.pdf | 479.011 | 47 | 60 | 491 | 13,2% | 10 | PDF 1.4 | text oriented, many links |
| sdw_2015_01.pdf | 14.668.972 | 100 | 36 | 0 | 2,5% | 143 | PDF 1.3 | graphics oriented |
| sdw_2015_02.pdf | 13.295.864 | 100 | 38 | 0 | 2,7% | 130 | PDF 1.4 | graphics oriented |
| sdw_2015_03.pdf | 21.224.417 | 108 | 35 | 0 | 1,9% | 192 | PDF 1.4 | graphics oriented |
| sdw_2015_04.pdf | 15.242.911 | 108 | 37 | 0 | 2,7% | 138 | PDF 1.3 | graphics oriented |
| sdw_2015_05.pdf | 16.495.887 | 108 | 43 | 0 | 2,4% | 149 | PDF 1.4 | graphics oriented |
| sdw_2015_06.pdf | 23.447.046 | 100 | 38 | 0 | 1,6% | 229 | PDF 1.4 | graphics oriented |
| sdw_2015_07.pdf | 14.106.982 | 100 | 38 | 2 | 2,6% | 138 | PDF 1.4 | graphics oriented |
| sdw_2015_08.pdf | 12.321.995 | 100 | 37 | 0 | 3,0% | 120 | PDF 1.4 | graphics oriented |
| sdw_2015_09.pdf | 23.409.625 | 100 | 37 | 0 | 1,5% | 229 | PDF 1.4 | graphics oriented |
| sdw_2015_10.pdf | 18.706.394 | 100 | 24 | 0 | 2,0% | 183 | PDF 1.5 | graphics oriented |
| sdw_2015_11.pdf | 25.624.266 | 100 | 20 | 0 | 1,5% | 250 | PDF 1.4 | graphics oriented |
| sdw_2015_12.pdf | 19.111.666 | 108 | 36 | 0 | 2,1% | 173 | PDF 1.4 | graphics oriented |

Decimal point and comma follow European convention

E.g. Adobe.pdf and PyMuPDF.pdf are clearly text oriented, all other files contain many more images.

8.1 Part 1: Parsing

How fast is a PDF file read and its content parsed for further processing? The sheer parsing performance cannot directly be compared, because batch utilities always execute a requested task completely, in one go, front to end. **pdfcrowd** too, has a **lazy** strategy for parsing, meaning it only parses those parts of a document that are required in any moment.

In order to yet find an answer to the question, we therefore measure the time to copy a PDF file to an output file with each tool, and doing nothing else.

These were the tools

All tools are either platform independent, or at least can run both, on Windows and Unix / Linux (pdftk).

Poppler is missing here, because it specifically is a Linux tool set, although we know there exist Windows ports (created with considerable effort apparently). Technically, it is a C/C++ library, for which a Python binding exists - in so far somewhat comparable to PyMuPDF. But Poppler in contrast is tightly coupled to **Qt** and **Cairo**. We may still include it in future, when a more handy Windows installation is available. We have seen however some [analysis](#), that hints at a much lower performance than MuPDF. Our comparison of text extraction speeds also show a much lower performance of Poppler's PDF code base **Xpdf**.

Image rendering of MuPDF also is about three times faster than the one of Xpdf when comparing the command line tools **mudraw** of MuPDF and **pdftopng** of Xpdf - see part 3 of this chapter.

| Tool | Description |
|---------|--|
| PyMuPDF | tool of this manual, appearing as "fitz" in reports |
| pdfrw | a pure Python tool, is being used by rst2pdf, has interface to ReportLab |
| PyPDF2 | a pure Python tool with a very complete function set |
| pdftk | a command line utility with numerous functions |

This is how each of the tools was used:

PyMuPDF:

```
doc = fitz.open("input.pdf")
doc.save("output.pdf")
```

pdfrw:

```
doc = PdfReader("input.pdf")
writer = PdfWriter()
writer.trailer = doc
writer.write("output.pdf")
```

PyPDF2:

```
pdfmerge = PyPDF2.PdfFileMerger()
pdfmerge.append("input.pdf")
pdfmerge.write("output.pdf")
pdfmerge.close()
```

pdftk:

```
pdftk input.pdf output output.pdf
```

Observations

These are our run time findings (in **seconds**, please note the European number convention: meaning of decimal point and comma is reversed):

| Runtime | Tool | | | |
|-----------------------|--------------|--------------|---------------|---------------|
| File | 1 fitz | 2 pdfrw | 3 pdftk | 4 PyPDF2 |
| Adobe.pdf | 5,25 | 21,06 | 112,39 | 692,23 |
| Evolution.pdf | 0,16 | 0,46 | 1,05 | 0,89 |
| PyMuPDF.pdf | 0,04 | 0,19 | 0,82 | 0,88 |
| sdw_2015_01.pdf | 0,23 | 1,23 | 5,41 | 6,45 |
| sdw_2015_02.pdf | 0,29 | 1,52 | 7,05 | 6,70 |
| sdw_2015_03.pdf | 0,51 | 2,77 | 11,49 | 11,98 |
| sdw_2015_04.pdf | 0,31 | 2,15 | 7,44 | 7,21 |
| sdw_2015_05.pdf | 0,35 | 1,69 | 7,60 | 7,59 |
| sdw_2015_06.pdf | 0,75 | 3,31 | 13,97 | 14,54 |
| sdw_2015_07.pdf | 0,37 | 2,11 | 10,17 | 9,72 |
| sdw_2015_08.pdf | 0,46 | 1,94 | 8,80 | 8,69 |
| sdw_2015_09.pdf | 0,79 | 2,35 | 10,58 | 10,42 |
| sdw_2015_10.pdf | 0,36 | 1,88 | 3,53 | 6,64 |
| sdw_2015_11.pdf | 2,41 | 12,69 | 37,12 | 60,40 |
| sdw_2015_12.pdf | 0,51 | 2,19 | 9,25 | 10,03 |
| Gesamtergebnis | 12,78 | 57,54 | 246,66 | 854,36 |

| | | | |
|------|------|-------|-------|
| 1,00 | 4,50 | 19,30 | 66,85 |
| | 1,00 | 4,29 | 14,85 |
| | | 1,00 | 3,46 |

If we leave out the Adobe manual, this table looks like

| Runtime | Tool | | | |
|-----------------------|-------------|--------------|---------------|---------------|
| File | 1 fitz | 2 pdfrw | 3 pdftk | 4 PyPDF2 |
| Evolution.pdf | 0,16 | 0,46 | 1,05 | 0,89 |
| PyMuPDF.pdf | 0,04 | 0,19 | 0,82 | 0,88 |
| sdw_2015_01.pdf | 0,23 | 1,23 | 5,41 | 6,45 |
| sdw_2015_02.pdf | 0,29 | 1,52 | 7,05 | 6,70 |
| sdw_2015_03.pdf | 0,51 | 2,77 | 11,49 | 11,98 |
| sdw_2015_04.pdf | 0,31 | 2,15 | 7,44 | 7,21 |
| sdw_2015_05.pdf | 0,35 | 1,69 | 7,60 | 7,59 |
| sdw_2015_06.pdf | 0,75 | 3,31 | 13,97 | 14,54 |
| sdw_2015_07.pdf | 0,37 | 2,11 | 10,17 | 9,72 |
| sdw_2015_08.pdf | 0,46 | 1,94 | 8,80 | 8,69 |
| sdw_2015_09.pdf | 0,79 | 2,35 | 10,58 | 10,42 |
| sdw_2015_10.pdf | 0,36 | 1,88 | 3,53 | 6,64 |
| sdw_2015_11.pdf | 2,41 | 12,69 | 37,12 | 60,40 |
| sdw_2015_12.pdf | 0,51 | 2,19 | 9,25 | 10,03 |
| Gesamtergebnis | 7,53 | 36,48 | 134,28 | 162,13 |

| | | | |
|------|------|-------|-------|
| 1,00 | 4,84 | 17,82 | 21,52 |
| | 1,00 | 3,68 | 4,44 |
| | | 1,00 | 1,21 |

PyMuPDF is by far the fastest: on average 4.5 times faster than the second best (the pure Python tool `pdfrw`, **chapeau pdfrw!**), and almost 20 times faster than the command line tool `pdftk`.

Where PyMuPDF only requires less than 13 seconds to process all files, `pdftk` affords itself almost 4 minutes.

By far the slowest tool is PyPDF2 - it is more than 66 times slower than PyMuPDF and 15 times slower than `pdfrw`! The main reason for PyPDF2's bad look comes from the Adobe manual. It obviously is slowed down by the linear file structure and the immense amount of bookmarks of this file. If we take out this special case, then PyPDF2 is only 21.5 times slower than PyMuPDF, 4.5 times slower than `pdfrw` and 1.2 times slower than `pdftk`.

If we look at the output PDFs, there is one surprise:

Each tool created a PDF of similar size as the original. Apart from the Adobe case, PyMuPDF always created the smallest output.

Adobe's manual is an exception: The pure Python tools `pdfrw` and PyPDF2 **reduced** its size by more than 20% (and yielded a document which is no longer linearized)!

PyMuPDF and `pdftk` in contrast **drastically increased** the size by 40% to about 50 MB (also no longer linearized).

So far, we have no explanation of what is happening here.

8.2 Part 2: Text Extraction

We also have compared text extraction speed with other tools.

The following table shows a run time comparison. PyMuPDF's methods appear as "fitz (TEXT)" and "fitz (JSON)" respectively. The tool `pdftotext.exe` of the [Xpdf](#) toolset appears as "xpdf".

- **extractText()**: basic text extraction without layout re-arrangement (using `GetText(..., output = "text")`)
- **pdftotext**: a command line tool of the **Xpdf** toolset (which also is the basis of [Poppler's library](#))
- **extractJSON()**: text extraction with layout information (using `GetText(..., output = "json")`)
- **pdfminer**: a pure Python PDF tool specialized on text extraction tasks

All tools have been used with their most basic, fanciless functionality - no layout re-arrangements, etc.

For demonstration purposes, we have included a version of `GetText(doc, output = "json")`, that also re-arranges the output according to occurrence on the page.

Here are the results using the same test files as above (again: decimal point and comma reversed):

| Runtime | Tool | | | | |
|-----------------------|---------------|-----------------|-----------------|--------------|----------------|
| File | 1 fitz (TEXT) | 2 fitz bareJSON | 3 fitz sortJSON | 4 xpdf | 5 pdfminer |
| Adobe.pdf | 5,16 | 5,53 | 6,27 | 12,42 | 216,32 |
| Evolution.pdf | 0,29 | 0,29 | 0,33 | 1,99 | 12,91 |
| PyMuPDF.pdf | 0,11 | 0,10 | 0,12 | 1,71 | 4,71 |
| sdw_2015_01.pdf | 0,95 | 0,98 | 1,12 | 2,84 | 43,96 |
| sdw_2015_02.pdf | 1,04 | 1,09 | 1,14 | 2,86 | 48,26 |
| sdw_2015_03.pdf | 1,81 | 1,92 | 1,97 | 3,82 | 153,51 |
| sdw_2015_04.pdf | 1,23 | 1,27 | 1,37 | 3,17 | 80,95 |
| sdw_2015_05.pdf | 1,00 | 1,08 | 1,15 | 2,82 | 48,65 |
| sdw_2015_06.pdf | 1,83 | 1,92 | 1,98 | 3,70 | 138,75 |
| sdw_2015_07.pdf | 0,99 | 1,11 | 1,16 | 2,93 | 55,59 |
| sdw_2015_08.pdf | 0,97 | 1,04 | 1,12 | 2,80 | 48,09 |
| sdw_2015_09.pdf | 1,92 | 1,97 | 2,05 | 3,84 | 159,62 |
| sdw_2015_10.pdf | 1,10 | 1,18 | 1,25 | 3,45 | 74,25 |
| sdw_2015_11.pdf | 2,37 | 2,39 | 2,50 | 5,82 | 166,14 |
| sdw_2015_12.pdf | 1,14 | 1,19 | 1,26 | 2,93 | 69,79 |
| Gesamtergebnis | 21,92 | 23,08 | 24,82 | 57,10 | 1321,51 |

| | | | | |
|------|------|------|------|-------|
| 1,00 | 1,05 | 1,13 | 2,60 | 60,28 |
| | 1,00 | 1,08 | 2,47 | 57,27 |
| | | 1,00 | 2,30 | 53,24 |
| | | | 1,00 | 23,15 |

Again, (Py-) MuPDF is the fastest around. It is 2.3 to 2.6 times faster than xpdf.

pdfminer, as a pure Python solution, of course is comparatively slow: MuPDF is 50 to 60 times faster and xpdf is 23 times faster. These observations in order of magnitude coincide with the statements on [this web site](#).

8.3 Part 3: Image Rendering

We have tested rendering speed of MuPDF against the `pdftopng.exe`, a command line tool of the **Xpdf** toolset (the PDF code basis of **Poppler**).

MuPDF invocation using a resolution of 150 pixels (Xpdf default):

```
mutool draw -o t%d.png -r 150 file.pdf
```

PyMuPDF invocation:

```
zoom = 150.0 / 72.0
mat = fitz.Matrix(zoom, zoom)
def ProcessFile(datei):
    print "processing:", datei
    doc=fitz.open(datei)
    for p in fitz.Pages(doc):
        pix = p.getPixmap(matrix=mat, alpha = False)
        pix.writePNG("t-%s.png" % p.number)
        pix = None
    doc.close()
    return
```

Xpdf invocation:

```
pdftopng.exe file.pdf ./
```

The resulting runtimes can be found here (again: meaning of decimal point and comma reversed):

| Render Speed | tool | | |
|-----------------------|---------------|---------------|----------------|
| file | mudraw | pymupdf | xpdf |
| Adobe.pdf | 105,09 | 110,66 | 505,27 |
| Evolution.pdf | 40,70 | 42,17 | 108,33 |
| PyMuPDF.pdf | 5,09 | 4,96 | 21,82 |
| sdw_2015_01.pdf | 29,77 | 30,40 | 76,81 |
| sdw_2015_02.pdf | 29,67 | 30,00 | 74,68 |
| sdw_2015_03.pdf | 32,67 | 32,88 | 85,89 |
| sdw_2015_04.pdf | 30,07 | 29,59 | 78,09 |
| sdw_2015_05.pdf | 31,37 | 31,39 | 77,56 |
| sdw_2015_06.pdf | 31,76 | 31,49 | 87,89 |
| sdw_2015_07.pdf | 33,33 | 34,58 | 78,74 |
| sdw_2015_08.pdf | 31,83 | 32,73 | 75,95 |
| sdw_2015_09.pdf | 36,92 | 36,77 | 84,37 |
| sdw_2015_10.pdf | 30,08 | 30,48 | 77,13 |
| sdw_2015_11.pdf | 33,21 | 34,11 | 80,96 |
| sdw_2015_12.pdf | 31,77 | 32,69 | 80,68 |
| Gesamtergebnis | 533,33 | 544,90 | 1594,18 |

| | | |
|---|------|------|
| 1 | 1,02 | 2,99 |
| | 1 | 2,93 |

- MuPDF and PyMuPDF are both about 3 times faster than Xpdf.
- The 2% speed difference between MuPDF (a utility written in C) and PyMuPDF is the Python overhead.

APPENDIX 2: DETAILS ON TEXT EXTRACTION

This chapter provides background on the text extraction methods of PyMuPDF.

Information of interest are

- what do they provide?
- what do they imply (processing time / data sizes)?

9.1 General structure of a *TextPage*

Text information contained in a *TextPage* adheres to the following hierarchy:

```
<page> (width and height)
  <block> (its rectangle)
    <line> (its rectangle)
      <span> (its rectangle and font information)
        <char> (its rectangle, (x, y) coordinates and value)
```

A **text page** consists of blocks (= roughly paragraphs).

A **block** consists of lines.

A **line** consists of spans.

A **span** consists of characters with the same properties. E.g. a different font will cause a new span.

9.2 Output of `getText(output="text")`

This function extracts a page's plain **text in original order** as specified by the creator of the document (which may not be equal to a natural reading order!).

An example output of this tutorial's PDF version:

```
Tutorial

This tutorial will show you the use of MuPDF in Python step by step.

Because MuPDF supports not only PDF, but also XPS, OpenXPS and EPUB formats, so does PyMuPDF.

Nevertheless we will only talk about PDF files for the sake of brevity.
...
```

9.3 Output of `getText(output="html")`

HTML output reflects the structure of the page's `TextPage` - without adding much other benefit. Again an example:

```
<div class="page">
<div class="block"><p>
<div class="metaline"><div class="line"><div class="cell" style="width:0%;align:left"><span
↪class="s0">Tutorial</span></div></div>
</div></p></div>
<div class="block"><p>
<div class="line"><div class="cell" style="width:0%;align:left"><span class="s1">This tutorial
↪will show you the use of MuPDF in Python step by step.</span></div></div>
</div></p></div>
<div class="block"><p>
<div class="line"><div class="cell" style="width:0%;align:left"><span class="s1">Because MuPDF
↪supports not only PDF, but also XPS, OpenXPS and EPUB formats, so does PyMuPDF.</span></div>
↪</div>
<div class="line"><div class="cell" style="width:0%;align:left"><span class="s1">Nevertheless
↪we will only talk about PDF files for the sake of brevity.</span></div></div>
</div></p></div>
...
```

9.4 Output of `getText(output="json")`

JSON output reflects the structure of a `TextPage` and provides position details (`bbox` - boundary boxes in pixel units) for every block, line and span. This is enough information to present a page's text in any required reading order (e.g. from top-left to bottom-right). The output can obviously be made usable by `text_dict = json.loads(text)`. Have a look at our example program `PDF2textJS.py`. Here is how it looks like:

```
{
  "len":35,"width":595.2756,"height":841.8898,
  "blocks":[
    {"type":"text","bbox":[40.01575, 53.730354, 98.68775, 76.08236],
      "lines":[
        {"bbox":[40.01575, 53.730354, 98.68775, 76.08236],
          "spans":[
            {"bbox":[40.01575, 53.730354, 98.68775, 76.08236],
              "text":"Tutorial"
            }
          ]
        }
      ]
    },
    {"type":"text","bbox":[40.01575, 79.300354, 340.6957, 93.04035],
      "lines":[
        {"bbox":[40.01575, 79.300354, 340.6957, 93.04035],
          "spans":[
            {"bbox":[40.01575, 79.300354, 340.6957, 93.04035],
              "text":"This tutorial will show you the use of MuPDF in Python step by step."
            }
          ]
        }
      ]
    }
  ],
  ...
}
```

9.5 Output of `getText(output="xml")`

The XML version takes the level of detail even a lot deeper: every single character is provided with its position detail, and every span also contains font information:

```
<page width="595.2756" height="841.8898">
<block bbox="40.01575 53.730354 98.68775 76.08236">
<line bbox="40.01575 53.730354 98.68775 76.08236">
<span bbox="40.01575 53.730354 98.68775 76.08236" font="Helvetica-Bold" size="16">
<char bbox="40.01575 53.730354 49.79175 76.08236" x="40.01575" y="70.85036" c="T"/>
<char bbox="49.79175 53.730354 59.56775 76.08236" x="49.79175" y="70.85036" c="u"/>
<char bbox="59.56775 53.730354 64.89575 76.08236" x="59.56775" y="70.85036" c="t"/>
<char bbox="64.89575 53.730354 74.67175 76.08236" x="64.89575" y="70.85036" c="o"/>
<char bbox="74.67175 53.730354 80.89575 76.08236" x="74.67175" y="70.85036" c="r"/>
<char bbox="80.89575 53.730354 85.34375 76.08236" x="80.89575" y="70.85036" c="i"/>
<char bbox="85.34375 53.730354 94.23975 76.08236" x="85.34375" y="70.85036" c="a"/>
<char bbox="94.23975 53.730354 98.68775 76.08236" x="94.23975" y="70.85036" c="l"/>
</span>
</line>
</block>
<block bbox="40.01575 79.300354 340.6957 93.04035">
<line bbox="40.01575 79.300354 340.6957 93.04035">
<span bbox="40.01575 79.300354 340.6957 93.04035" font="Helvetica" size="10">
<char bbox="40.01575 79.300354 46.12575 93.04035" x="40.01575" y="90.050354" c="T"/>
<char bbox="46.12575 79.300354 51.685753 93.04035" x="46.12575" y="90.050354" c="h"/>
<char bbox="51.685753 79.300354 53.90575 93.04035" x="51.685753" y="90.050354" c="i"/>
<char bbox="53.90575 79.300354 58.90575 93.04035" x="53.90575" y="90.050354" c="s"/>
<char bbox="58.90575 79.300354 61.685753 93.04035" x="58.90575" y="90.050354" c=" " />
<char bbox="61.685753 79.300354 64.46575 93.04035" x="61.685753" y="90.050354" c="t"/>
<char bbox="64.46575 79.300354 70.02576 93.04035" x="64.46575" y="90.050354" c="u"/>
<char bbox="70.02576 79.300354 72.805756 93.04035" x="70.02576" y="90.050354" c="t"/>
<char bbox="72.805756 79.300354 78.36575 93.04035" x="72.805756" y="90.050354" c="o"/>
<char bbox="78.36575 79.300354 81.695755 93.04035" x="78.36575" y="90.050354" c="r"/>
<char bbox="81.695755 79.300354 83.91576 93.04035" x="81.695755" y="90.050354" c="i"/>
...

```

The method's output can be processed by one of Python's XML modules. We have successfully tested `lxml`. See the demo program `fontlistner.py`. It creates a list of all fonts of a document including font size and where used on pages.

9.6 Performance

The four text extraction methods of a *TextPage* differ significantly: in terms of information they supply (see above), and in terms of resource requirements. More information of course means that more processing is required and a higher data volume is generated.

To begin with, all four methods are **very** fast in relation to what is there on the market. In terms of processing speed, we couldn't find a faster (free) tool.

Relative to each other, `xml` is about 2 times slower than `text`, the other three range between them. E.g. `json` needs about 13% - 14% more time than `text`.

Look into the previous chapter **Appendix 1** for more performance information.

APPENDIX 3: CONSIDERATIONS ON EMBEDDED FILES

This chapter provides some background on embedded files support in PyMuPDF.

10.1 General

Starting with version 1.4, PDF supports embedding arbitrary files as part (“Embedded File Streams”) of a PDF document file (see chapter 3.10.3, pp. 184 of the *Adobe PDF Reference 1.7*).

In many aspects, this is comparable to concepts also found in ZIP files or the OLE technique in MS Windows. PDF embedded files do, however, *not* support directory structures as does the ZIP format. An embedded file can in turn contain embedded files itself.

Advantages of this concept are that embedded files are under the PDF umbrella, benefitting from its permissions / password protection and integrity aspects: all files a PDF may reference or even be dependent on can be bundled into it and so form a single, consistent unit of information.

In addition to embedded files, PDF 1.7 adds *collections* to its support range. This is an advanced way of storing and presenting meta information (i.e. arbitrary and extensible properties) of embedded files.

10.2 MuPDF Support

MuPDF v1.11 added initial support for embedded files and collections (also called *portfolios*).

The library contains functions to add files to the `EmbeddedFiles` name tree and display some information of its entries.

Also supported is a full set of functions to maintain collections (advanced metadata maintenance) and their relation to embedded files.

10.3 PyMuPDF Support

PyMuPDF v1.11.0 fully reflects MuPDF’s support for embedded files and partly goes beyond that scope:

- We can add, extract **and** delete embedded files.
- We can display **and** change some meta information (outside collections). Informations available for display are **name**, **filename**, **description**, **length** and compressed **size**. Of these properties, *filename* and *description* can also be changed, after a file has been embedded.

Support of the *collections* feature has been postponed to a later version. We will probably include this ever only on user request.

APPENDIX 4: ASSORTED TECHNICAL INFORMATION

11.1 PDF Base 14 Fonts

The following 14 builtin font names must be supported by every PDF application. They are available as the Python list `fitz.Base14_Fonts`:

- Courier
- Courier-Oblique
- Courier-Bold
- Courier-BoldOblique
- Helvetica
- Helvetica-Oblique
- Helvetica-Bold
- Helvetica-BoldOblique
- Times-Roman
- Times-Bold
- Times-Italic
- Times-BoldItalic
- Symbol
- ZapfDingbats

11.2 Adobe PDF Reference 1.7

This PDF Reference manual published by Adobe is frequently quoted throughout this documentation. It can be viewed and downloaded from here: http://www.adobe.com/content/dam/Adobe/en/devnet/acrobat/pdfs/pdf_reference_1-7.pdf.

CHANGE LOGS

12.1 Changes in Version 1.11.0

This version is based on and requires MuPDF v1.11.

Though MuPDF has declared it as being mostly a bug fix version, one major new feature is indeed contained: support of embedded files - also called portfolios or collections. We have extended PyMuPDF functionality to embrace this up to an extent just a little beyond the `mutool` utility as follows.

- The `Document` class contains several new methods and one new property to support portfolios:
 - `embeddedFileInfo()` returns metadata information about an entry in the list of embedded files. This is more than `mutool` currently provides: it shows all the information that was used to embed the file (not just the entry's name).
 - `embeddedFileGet()` retrieves the (decompressed) content of an entry into a `bytes` buffer.
 - `embeddedFileAdd(...)` inserts new content into the PDF portfolio. We (in contrast to `mutool`) **restrict** this to entries with a **new name** (no duplicate names allowed).
 - `embeddedFileDel(...)` deletes an entry from the portfolio (function not offered in MuPDF).
 - `embeddedFileSetInfo()` - changes filename or description of an embedded file.
 - `embeddedFileCount` - contains the number of embedded files.
- Several enhancements deal with streamlining geometry objects. These are not connected to the new MuPDF version and most of them are also reflected in PyMuPDF v1.10.0. Among them are new properties to identify the corners of rectangles by name (e.g. `Rect.bottom_right`) and new methods to deal with set-theoretic questions like `Rect.contains(x)` or `IRect.intersects(x)`. Special effort focussed on supporting more “Pythonic” language constructs: `if x in rect ...` is equivalent to `rect.contains(x)`.
- The [Rect](#) chapter now has more background on empty and infinite rectangles and how we handle them. The handling itself was also updated for more consistency in this area.
- We have started basic support for **generation** of PDF content:
 - `Document.insertPage()` adds a new page into a PDF, optionally containing some text.
 - `Page.insertImage()` places a new image on a PDF page.
 - `Page.insertText()` puts new text on an existing page
- For **FileAttachment** annotations, content and name of the attached file can be extracted and changed.

12.2 Changes in Version 1.10.0

12.2.1 MuPDF v1.10 Impact

MuPDF version 1.10 has a significant impact on our bindings. Some of the changes also affect the API - in other words, **you** as a PyMuPDF user.

- Link destination information has been reduced. Several properties of the `linkDest` class no longer contain valuable information. In fact, this class as a whole has been deleted from MuPDF's library and we in PyMuPDF only maintain it to provide compatibility to existing code.
- In an effort to minimize memory requirements, several improvements have been built into MuPDF v1.10:
 - A new `config.h` file can be used to de-select unwanted features in the C base code. Using this feature we have been able to reduce the size of our binary `_fitz.o` / `_fitz.pyd` by about 50% (from 9 MB to 4.5 MB). When UPX-ing this, the size goes even further down to a very handy 2.3 MB.
 - The alpha (transparency) channel for pixmaps is now optional. Letting alpha default to `False` significantly reduces pixmap sizes (by 20% - CMYK, 25% - RGB, 50% - GRAY). Many `Pixmap` constructors therefore now accept an `alpha` boolean to control inclusion of this channel. Other pixmap constructors (e.g. those for file and image input) create pixmaps with no alpha altogether. On the downside, save methods for pixmaps no longer accept a `savealpha` option: this channel will always be saved when present. In order to minimize code breaks, we have left this parameter in the call patterns - it will just be ignored.
- `DisplayList` and `TextPage` class constructors now **require the mediabox** of the page they are referring to (i.e. the `page.bound()` rectangle). There is no way to construct this information from other sources, therefore a source code change cannot be avoided in these cases. We assume however, that not many users are actually employing these rather low level classes explicitly. So the impact of that change should be minor.

12.2.2 Other Changes compared to Version 1.9.3

- The new `Document` method `write()` writes an opened PDF to memory (as opposed to a file, like `save()` does).
- An annotation can now be scaled and moved around on its page. This is done by modifying its rectangle.
- Annotations can now be deleted. `Page` contains the new method `deleteAnnot()`.
- Various annotation attributes can now be modified, e.g. content, dates, title (= author), border, colors.
- Method `Document.insertPDF()` now also copies annotations of source pages.
- The `Pages` class has been deleted. As documents can now be accessed with page numbers as indices (like `doc[n] = doc.loadPage(n)`), and document object can be used as iterators, the benefit of this class was too low to maintain it. See the following comments.
- `loadPage(n)` / `doc[n]` now accept arbitrary integers to specify a page number, as long as `n < pageCount`. So, e.g. `doc[-500]` is always valid and will load page `(-500) % pageCount`.
- A document can now also be used as an iterator like this: `for page in doc: ...<do something with "page">` This will yield all pages of `doc` as `page`.
- The `Pixmap` method `getSize()` has been replaced with property `size`. As before `Pixmap.size == len(Pixmap)` is `True`.
- In response to transparency (alpha) being optional, several new parameters and properties have been added to `Pixmap` and `Colorspace` classes to support determining their characteristics.

- The *Page* class now contains new properties `firstAnnot` and `firstLink` to provide starting points to the respective class chains, where `firstLink` is just a mnemonic synonym to method `loadLinks()` which continues to exist. Similarly, the new property `rect` is a synonym for method `bound()`, which also continues to exist.
- *Pixmap* methods `samplesRGB()` and `samplesAlpha()` have been deleted because pixmaps can now be created without transparency.
- *Rect* now has a property `irect` which is a synonym of method `round()`. Likewise, *IRect* now has property `rect` to deliver a *Rect* which has the same coordinates as floats values.
- Document has the new method `searchPageFor()` to search for a text string. It works exactly like the corresponding `Page.searchFor()` with page number as additional parameter.

12.3 Changes in Version 1.9.3

This version is also based on MuPDF v1.9a. Changes compared to version 1.9.2:

- As a major enhancement, annotations are now supported in a similar way as links. Annotations can be displayed (as pixmaps) and their properties can be accessed.
- In addition to the document `select()` method, some simpler methods can now be used to manipulate a PDF:
 - `copyPage()` copies a page within a document.
 - `movePage()` is similar, but deletes the original.
 - `deletePage()` deletes a page
 - `deletePageRange()` deletes a page range
- `rotation` or `setRotation()` access or change a PDF page's rotation, respectively.
- Available but undocumented before, *IRect*, *Rect*, *Point* and *Matrix* support the `len()` method and their coordinate properties can be accessed via indices, e.g. `IRect.x1 == IRect[2]`.
- For convenience, documents now support simple indexing: `doc.loadPage(n) == doc[n]`. The index may however be in range `-pageCount < n < pageCount`, such that `doc[-1]` is the last page of the document.

12.4 Changes in Version 1.9.2

This version is also based on MuPDF v1.9a. Changes compared to version 1.9.1:

- `fitz.open()` (no parameters) creates a new empty **PDF** document, i.e. if saved afterwards, it must be given a `.pdf` extension.
- *Document* now accepts all of the following formats (*Document* and *open* are synonyms):
 - `open()`,
 - `open(filename)` (equivalent to `open(filename, None)`),
 - `open filetype, area)` (equivalent to `open filetype, stream = area)`).

Type of memory area `stream` may be `str` (Python 2), `bytes` (Python 3) or `bytearray` (Python 2 and 3). Thus, e.g. `area = open("file.pdf", "rb").read()` may be used directly (without first converting it to `bytearray`).

- New method `Document.insertPDF()` (PDFs only) inserts a range of pages from another PDF.
- Document objects `doc` now support the `len()` function: `len(doc) == doc.pageCount`.
- New method `Document.getPageImageList()` creates a list of images used on a page.

- New method `Document.getPageFontList()` creates a list of fonts referenced by a page.
- New pixmap constructor `fitz.Pixmap(doc, xref)` creates a pixmap based on an opened PDF document and an XREF number of the image.
- New pixmap constructor `fitz.Pixmap(cspace, spix)` creates a pixmap as a copy of another one `spix` with the colorspace converted to `cspace`. This works for all colorspace combinations.
- Pixmap constructor `fitz.Pixmap(colorspace, width, height, samples)` now allows `samples` to also be `str` (Python 2) or `bytes` (Python 3), not only `bytearray`.

12.5 Changes in Version 1.9.1

This version of PyMuPDF is based on MuPDF library source code version 1.9a published on April 21, 2016.

Please have a look at MuPDF's website to see which changes and enhancements are contained herein.

Changes in version 1.9.1 compared to version 1.8.0 are the following:

- New methods `getRectArea()` for both `fitz.Rect` and `fitz.IRect`
- Pixmap can now be created directly from files using the new constructor `fitz.Pixmap(filename)`.
- The Pixmap constructor `fitz.Pixmap(image)` has been extended accordingly.
- `fitz.Rect` can now be created with all possible combinations of points and coordinates.
- PyMuPDF classes and methods now all contain `__doc__` strings, most of them created by SWIG automatically. While the PyMuPDF documentation certainly is more detailed, this feature should help a lot when programming in Python-aware IDEs.
- A new document method of `getPermits()` returns the permissions associated with the current access to the document (print, edit, annotate, copy), as a Python dictionary.
- The identity matrix `fitz.Identity` is now **immutable**.
- The new document method `select(list)` removes all pages from a document that are not contained in the list. Pages can also be duplicated and re-arranged.
- Various improvements and new members in our demo and examples collections. Perhaps most prominently: `PDF_display` now supports scrolling with the mouse wheel, and there is a new example program `wxTableExtract` which allows to graphically identify and extract table data in documents.
- `fitz.open()` is now an alias of `fitz.Document()`.
- New pixmap method `getPNGData()` which will return a bytearray formatted as a PNG image of the pixmap.
- New pixmap method `samplesRGB()` providing a `samples` version with alpha bytes stripped off (RGB colorspace only).
- New pixmap method `samplesAlpha()` providing the alpha bytes only of the `samples` area.
- New iterator `fitz.Pages(doc)` over a document's set of pages.
- New matrix methods `invert()` (calculate inverted matrix), `concat()` (calculate matrix product), `preTranslate()` (perform a shift operation).
- New `IRect` methods `intersect()` (intersection with another rectangle), `translate()` (perform a shift operation).
- New `Rect` methods `intersect()` (intersection with another rectangle), `transform()` (transformation with a matrix), `includePoint()` (enlarge rectangle to also contain a point), `includeRect()` (enlarge rectangle to also contain another one).
- Documented `Point.transform()` (transform a point with a matrix).

- `Matrix`, `IRect`, `Rect` and `Point` classes now support compact, algebraic formulations for manipulating such objects.
- Incremental saves for changes are possible now using the call pattern `doc.save(doc.name, incremental=True)`.
- A PDF's metadata can now be deleted, set or changed by document method `setMetadata()`. Supports incremental saves.
- A PDF's bookmarks (or table of contents) can now be deleted, set or changed with the entries of a list using document method `setToC(list)`. Supports incremental saves.

ERROR MESSAGES

This is a list of exception messages raised by PyMuPDF together with an explanation and possible solution. In addition, the underlying C library MuPDF also raises exceptions on the Python level. We have included a few of those as well and may extend this in future.

annot has no /AP

- Bad specification - no changes possible for this annotation.

arg 1 not bytes or bytearray

- Specify parameter as type `bytes` or `bytearray`.

bad PDF: Contents is no stream object

- The `/Contents` object(s) of a page must be streams. Repair PDF.

bad PDF: file has no stream

- An embedded / attached file is not a stream. Repair PDF.

buffer too large to deflate

- Internal error - report an issue.

cannot deflate buffer

- Internal error - report an issue.

cannot open <path>: No such file or directory

- Specify a valid file name / path.

cannot recognize archive

- Trying to open an invalid CBZ document.

cannot recognize zip archive

- Trying to open an invalid XPS document.

color components must be in range 0 to 1

- Color components must be floats in interval $[0, 1]$.

could not create UTF16 for '<name>'

- Internal error - report an issue.

could not get string of '<name>'

- Internal error - report an issue.

could not get UTF16 string of '<name>'

- Internal error - report an issue.

could not load root object

- Root object of PDF not found. Repair PDF.

encrypted file - save to new

- Trying incremental save for a decrypted file. Use `doc.save()` to a new file.

exactly one of filename, pixmap must be given

- You either specified both parameters or none.

expected a sequence

- Parameter type must be `list`, `tuple`, etc.

filename must be a string

- Specify a valid file path / name.

filename must be string or None

- Specify a valid file path / name or omit parameter.

filename must end with '.png'

- `writePNG()` requires file extension `.png`.

filetype missing with stream specified

- Document open from memory needs its type as a string.

fontname must be supplied

- A new font file requires some (arbitrary) **new** reference name.

found code point nnn: increase charlimit

- Trying to get a glyph width beyond the current table size limit.

incremental excludes garbage

- Garbage collection cannot occur during incremental saves.

incremental excludes linear

- Linearization cannot occur during incremental saves.

incremental save needs original file

- Incremental save is only possible to the original file.

info not a dict

- Specify correct Python parameter type.

invalid font - FontDescriptor missing

- Specify correct XREF to read font.

invalid font descriptor subtype

- Bad font description in PDF. Repair file.

unhandled font type / unhandled font type '<type>'

- MuPDF does not yet handle this font type. Requesting method cannot be used, unfortunately. Report an issue.

invalid key in info dict

- Dictionary key misspelled.

invalid page range

- Page numbers must be in range `[0, pageCount - 1]`.

invalid stream

- Stream object updates need type `bytes` or `bytearray`.

len(samples) invalid

- Length of samples must equal `width * height * n` (where `n` is the number of components per pixel).

line endpoints must be within page rect

- The `Page.rect` must contain the points.

name already exists

- The name is in use by some other embedded file.

name not valid

- Specify a name of non-zero length.

need 3 color components

- Only RGB colors are supported, which need three components.

no embedded files

- PDF has no embedded files.

no objects found

- Trying to open an invalid PDF, FB2, or EPUB document.

not a file attachment annot

- Accessed an annotation with the wrong type.

not a PDF

- Using some method or attribute only valid for PDF document type.

nothing to change

- No data supplied for embedded file metadata change.

operation illegal for closed doc

- Trying to use methods / properties after close of document.

orphaned object: parent is None

- Accessing an object whose parent no longer exists (e.g. an annotation of an unavailable page).

page number out of range

- Page numbers must always be `< pageCount`, but also non-negative for some methods.

page numbers must be integers

- Specify valid page numbers (`select()` method).

rect must be contained in page rect

- Image insertion requires a target rectangle contained in `page.rect`.

rect must be finite and not empty

- Top-left corner must be “northeast” of bottom-right one, and rectangle area must be positive.

repaired file - save to new

- Trying incremental save for file repaired during open. Use `doc.save()` to a new file.

save to original requires incremental

- Using original filename in `doc.save()` without also specifying option `incremental`. Consider using `doc.saveIncr()`.

sequence length must be `<n>`

- Creating Point, Rect, Irect, Matrix with wrong length sequences.

some text is needed

- Specify text with a positive length.

source and target too close

- Target number of moved page `pno` must be `> pno` or `< pno - 1`.

source must not equal target PDF

- Method `doc.insertPDF()` requires two distinct document objects (which may point to the same file, however).

source not a PDF

- Method `doc.insertPDF()` only works with PDF documents.

source page out of range

- Specify a valid page number.

target not a PDF

- Method `doc.insertPDF()` only works with PDF documents.

text position outside page height range

- If text starts at *Point* point, `fontsize <= point.y <= (page height - fontsize * 1.2)` must be true.

type(ap) invalid

- Internal error - report an issue.

type(imagedata) invalid

- Use type `bytearray`.

type(samples) invalid

- Use type `bytes` or `bytearray`.

unknown PDF Base 14 font

- Use a valid PDF standard font name.

xref entry is not an image

- Trying to create a pixmap from a non-image PDF object.

xref invalid

- Internal error - report an issue.

xref is not a stream

- Trying to access the stream part of a non-stream object.

xref out of range

- PDF xref numbers must be `1 <= xref <= doc._getXrefLength()`.

INDEX

- `__init__()` (Colorspace method), 17
- `__init__()` (Device method), 80
- `__init__()` (DisplayList method), 80
- `__init__()` (Document method), 19
- `__init__()` (IRect method), 33
- `__init__()` (Matrix method), 40
- `__init__()` (Pixmap method), 60, 61
- `__init__()` (Point method), 67
- `__init__()` (Rect method), 70
- `cleanContents()` (Annot method), 77
- `cleanContents()` (Page method), 77
- `delXmlMetadata()` (Document method), 76
- `getCharWidths()` (Document method), 77
- `getContents()` (Page method), 76
- `getNewXref()` (Document method), 78
- `getOLRootNumber()` (Document method), 79
- `getObjectString()` (Document method), 78
- `getPageObjNumber()` (Document method), 76
- `getPageRectText()` (Document method), 78
- `getPageXref()` (Document method), 76
- `getRectText()` (Page method), 78
- `getXref()` (Annot method), 77
- `getXref()` (Page method), 76
- `getXrefLength()` (Document method), 79
- `getXrefStream()` (Document method), 79
- `getXrefString()` (Document method), 78
- `updateObject()` (Document method), 78
- `updateStream()` (Document method), 79
- a (Matrix attribute), 41
- alpha (Pixmap attribute), 63
- Annot (built-in class), 13
- ANNOT_3D (built-in variable), 88
- ANNOT_CARET (built-in variable), 88
- ANNOT_CIRCLE (built-in variable), 87
- ANNOT_FILEATTACHMENT (built-in variable), 88
- ANNOT_FREETEXT (built-in variable), 87
- ANNOT_HIGHLIGHT (built-in variable), 87
- ANNOT_INK (built-in variable), 88
- ANNOT_LE_Butt (built-in variable), 90
- ANNOT_LE_Circle (built-in variable), 89
- ANNOT_LE_ClosedArrow (built-in variable), 89
- ANNOT_LE_Diamond (built-in variable), 89
- ANNOT_LE_None (built-in variable), 89
- ANNOT_LE_OpenArrow (built-in variable), 89
- ANNOT_LE_RClosedArrow (built-in variable), 90
- ANNOT_LE_ROpenArrow (built-in variable), 90
- ANNOT_LE_Slash (built-in variable), 90
- ANNOT_LE_Square (built-in variable), 89
- ANNOT_LINE (built-in variable), 87
- ANNOT_LINK (built-in variable), 87
- ANNOT_MOVIE (built-in variable), 88
- ANNOT_POLYGON (built-in variable), 87
- ANNOT_POLYLINE (built-in variable), 87
- ANNOT_POPUP (built-in variable), 88
- ANNOT_PRINTERMARK (built-in variable), 88
- ANNOT_SCREEN (built-in variable), 88
- ANNOT_SOUND (built-in variable), 88
- ANNOT_SQUARE (built-in variable), 87
- ANNOT_SQUIGGLY (built-in variable), 87
- ANNOT_STAMP (built-in variable), 88
- ANNOT_STRIKEOUT (built-in variable), 88
- ANNOT_TEXT (built-in variable), 87
- ANNOT_TRAPNET (built-in variable), 88
- ANNOT_UNDERLINE (built-in variable), 87
- ANNOT_WATERMARK (built-in variable), 88
- ANNOT_WIDGET (built-in variable), 88
- ANNOT_XF_Hidden (built-in variable), 88
- ANNOT_XF_Invisible (built-in variable), 88
- ANNOT_XF_Locked (built-in variable), 89
- ANNOT_XF_LockedContents (built-in variable), 89
- ANNOT_XF_NoRotate (built-in variable), 89
- ANNOT_XF_NoView (built-in variable), 89
- ANNOT_XF_NoZoom (built-in variable), 89
- ANNOT_XF_Print (built-in variable), 89
- ANNOT_XF_ReadOnly (built-in variable), 89
- ANNOT_XF_ToggleNoView (built-in variable), 89
- authenticate() (Document method), 19
- b (Matrix attribute), 41
- Base14_Fonts (built-in variable), 85
- border (Annot attribute), 16
- bottom_left (IRect attribute), 34
- bottom_left (Rect attribute), 72
- bottom_right (IRect attribute), 34
- bottom_right (Rect attribute), 72
- bound() (Page method), 49
- c (Matrix attribute), 41

clearWith() (Pixmap method), 61
close() (Document method), 27
colors (Annot attribute), 16
Colorspace (built-in class), 17
colorspace (Pixmap attribute), 63
concat() (Matrix method), 41
contains() (IRect method), 34
contains() (Rect method), 71
copyPage() (Document method), 26
copyPixmap() (Pixmap method), 62
CS_CMYK (built-in variable), 85
CS_GRAY (built-in variable), 85
CS_RGB (built-in variable), 85
csCMYK (built-in variable), 85
csGRAY (built-in variable), 85
csRGB (built-in variable), 85

d (Matrix attribute), 41
deleteAnnot() (Page method), 49
deleteLink() (Page method), 49
deletePage() (Document method), 25
deletePageRange() (Document method), 25
dest (Link attribute), 37
dest (linkDest attribute), 38
dest (Outline attribute), 48
Device (built-in class), 80
DisplayList (built-in class), 80
distance_to() (Point method), 67
Document (built-in class), 19
down (Outline attribute), 47
drawBezier() (Page method), 53
drawCircle() (Page method), 51
drawCurve() (Page method), 53
drawLine() (Page method), 51
drawOval() (Page method), 52
drawPolyline() (Page method), 52
drawRect() (Page method), 54
drawSector() (Page method), 52

e (Matrix attribute), 42
embeddedFileAdd() (Document method), 27
embeddedFileCount (Document attribute), 29
embeddedFileDel() (Document method), 27
embeddedFileGet() (Document method), 26
embeddedFileInfo() (Document method), 26
embeddedFileSetInfo() (Document method), 26
extractHTML() (TextPage method), 81
extractJSON() (TextPage method), 81
extractText() (TextPage method), 81
extractXML() (TextPage method), 81

f (Matrix attribute), 42
fileGet() (Annot method), 14
fileInfo() (Annot method), 14
fileSpec (linkDest attribute), 38
fileUpd() (Annot method), 15
firstAnnot (Page attribute), 56
firstLink (Page attribute), 56
flags (Annot attribute), 15

flags (linkDest attribute), 38

gammaWith() (Pixmap method), 62
getArea() (IRect method), 33
getArea() (Rect method), 71
getFontList() (Page method), 55
getImageList() (Page method), 55
getLinks() (Page method), 50
getPageFontList() (Document method), 21
getPageImageList() (Document method), 20
getPagePixmap() (Document method), 20
getPageText() (Document method), 21
getPDFnow(), 76
getPDFstr(), 76
getPixmap() (Annot method), 13
getPixmap() (Page method), 55
getPNGData() (Pixmap method), 62
getRect() (IRect method), 33
getRectArea() (IRect method), 33
getRectArea() (Rect method), 71
getText() (Page method), 55
getToC() (Document method), 20

height (IRect attribute), 34
height (Pixmap attribute), 63
height (Rect attribute), 72

includePoint() (Rect method), 71
includeRect() (Rect method), 71
info (Annot attribute), 15
insertImage() (Page method), 54
insertLink() (Page method), 50
insertPage() (Document method), 24
insertPDF() (Document method), 24
insertText() (Page method), 50
insertTextbox() (Page method), 50
interpolate (Pixmap attribute), 64
intersect() (IRect method), 33
intersect() (Rect method), 71
intersects() (IRect method), 34
intersects() (Rect method), 71
invert() (Matrix method), 41
invertIRect() (Pixmap method), 62
IRect (built-in class), 33
irect (Pixmap attribute), 63
is_open (Outline attribute), 47
isClosed (Document attribute), 27
isEmpty (IRect attribute), 35
isEmpty (Rect attribute), 72
isEncrypted (Document attribute), 28
isExternal (Link attribute), 37
isExternal (Outline attribute), 48
isInfinite (IRect attribute), 35
isInfinite (Rect attribute), 72
isMap (linkDest attribute), 38
isUri (linkDest attribute), 38

kind (linkDest attribute), 38

-
- lineEnds (Annot attribute), 16
 - Link (built-in class), 37
 - LINK_FLAG_B_VALID (built-in variable), 87
 - LINK_FLAG_FIT_H (built-in variable), 87
 - LINK_FLAG_FIT_V (built-in variable), 87
 - LINK_FLAG_L_VALID (built-in variable), 86
 - LINK_FLAG_R_IS_ZOOM (built-in variable), 87
 - LINK_FLAG_R_VALID (built-in variable), 87
 - LINK_FLAG_T_VALID (built-in variable), 87
 - LINK_GOTO (built-in variable), 86
 - LINK_GOTOR (built-in variable), 86
 - LINK_LAUNCH (built-in variable), 86
 - LINK_NONE (built-in variable), 86
 - LINK_URI (built-in variable), 86
 - linkDest (built-in class), 38
 - loadPage() (Document method), 19
 - lt (linkDest attribute), 38
 - Matrix (built-in class), 40
 - metadata (Document attribute), 28
 - movePage() (Document method), 26
 - n (Colorspace attribute), 18
 - n (Pixmap attribute), 63
 - name (Colorspace attribute), 18
 - name (Document attribute), 28
 - named (linkDest attribute), 38
 - needsPass (Document attribute), 27
 - newWindow (linkDest attribute), 39
 - next (Annot attribute), 15
 - next (Link attribute), 37
 - next (Outline attribute), 47
 - normalize() (IRect method), 34
 - normalize() (Rect method), 72
 - number (Page attribute), 56
 - openErrCode (Document attribute), 28
 - openErrMsg (Document attribute), 29
 - Outline (built-in class), 47
 - outline (Document attribute), 27
 - Page (built-in class), 49
 - page (linkDest attribute), 39
 - page (Outline attribute), 47
 - pageCount (Document attribute), 28
 - PaperSize(), 75
 - parent (Annot attribute), 15
 - parent (Page attribute), 56
 - permissions (Document attribute), 28
 - Pixmap (built-in class), 60
 - Point (built-in class), 67
 - preRotate() (Matrix method), 40
 - preScale() (Matrix method), 40
 - preShear() (Matrix method), 41
 - preTranslate() (Matrix method), 41
 - rb (linkDest attribute), 39
 - rect (Annot attribute), 15
 - Rect (built-in class), 70
 - rect (Link attribute), 37
 - rect (Page attribute), 49
 - rotation (Page attribute), 56
 - round() (Rect method), 70
 - run() (DisplayList method), 80
 - run() (Page method), 76
 - samples (Pixmap attribute), 63
 - save() (Document method), 23
 - saveIncr() (Document method), 23
 - saveText() (Outline method), 47
 - saveXML() (Outline method), 47
 - search() (TextPage method), 81
 - searchFor() (Page method), 56
 - searchPageFor() (Document method), 24
 - select() (Document method), 22
 - setBorder() (Annot method), 14
 - setColors() (Annot method), 14
 - setFlags() (Annot method), 14
 - setInfo() (Annot method), 14
 - setMetadata() (Document method), 22
 - setRect() (Annot method), 14
 - setRotation() (Page method), 55
 - setToC() (Document method), 22
 - size (Pixmap attribute), 63
 - stride (Pixmap attribute), 63
 - TEXT_ALIGN_CENTER (built-in variable), 86
 - TEXT_ALIGN_JUSTIFY (built-in variable), 86
 - TEXT_ALIGN_LEFT (built-in variable), 86
 - TEXT_ALIGN_RIGHT (built-in variable), 86
 - TextPage (built-in class), 81
 - tintWith() (Pixmap method), 61
 - title (Outline attribute), 47
 - top_left (IRect attribute), 34
 - top_left (Rect attribute), 72
 - top_right (IRect attribute), 34
 - top_right (Rect attribute), 72
 - transform() (Point method), 68
 - transform() (Rect method), 70
 - translate() (IRect method), 34
 - type (Annot attribute), 15
 - updateImage() (Annot method), 14
 - updateLink() (Page method), 50
 - uri (Link attribute), 37
 - uri (linkDest attribute), 39
 - uri (Outline attribute), 48
 - VersionBind (built-in variable), 85
 - VersionDate (built-in variable), 85
 - VersionFitz (built-in variable), 85
 - vertices (Annot attribute), 16
 - width (IRect attribute), 34
 - width (Pixmap attribute), 63
 - width (Rect attribute), 72
 - write() (Document method), 24

`writeImage()` (Pixmap method), 62
`writePNG()` (Pixmap method), 62

`x` (Pixmap attribute), 63
`x0` (IRect attribute), 34
`x0` (Rect attribute), 72
`x1` (IRect attribute), 35
`x1` (Rect attribute), 72
`xres` (Pixmap attribute), 64

`y` (Pixmap attribute), 63
`y0` (IRect attribute), 35
`y0` (Rect attribute), 72
`y1` (IRect attribute), 35
`y1` (Rect attribute), 72
`yres` (Pixmap attribute), 64