

ApacheCon

# Apache PDFBox - Working with pdfs for Dummies

Andreas Lehmkuhler



Leading the Wave  
of Open Source

## The Speaker

- Andreas Lehmkuhler
- Degree in electrical engineering
- Senior Developer
- Committer since December 2008
- PMC Chair since November 2009
- ASF member since April 2010



## Agenda

- History
- Portable Document Format
- Key features
- Using PDFBox inside applications
- Future prospects



## History of Apache PDFBox



## Yet another library dealing with pdfs?

- serves only one purpose
  - rendering
  - creating
  - editing
- can't be easily embedded in your own application
- wrong language
- license not suitable



## The Beginning

- Ben Litchfield started development in 2002
- initial purpose: extract text content to be indexed by the Lucene search engine
- hosted on Sourceforge
- BSD license
- Daniel Wilson and Philip Koch joined PDFBox



## Becoming an ASF Project

- several ASF projects showed their interest in a pdf-library
- software granted
- entering Apache Incubator in 2008
- first incubation release 0.8 in May 2009
- graduation as top level project in November 2009
- 10 committers
- current release 1.3.1



## Portable Document Format





## Portable Document Format (1)

- created by Adobe 1993
- open standard for document interchange
- ISO/IEC 32000-1:2008
- uses a subset of postscript
- includes used fonts
- single pages are independent

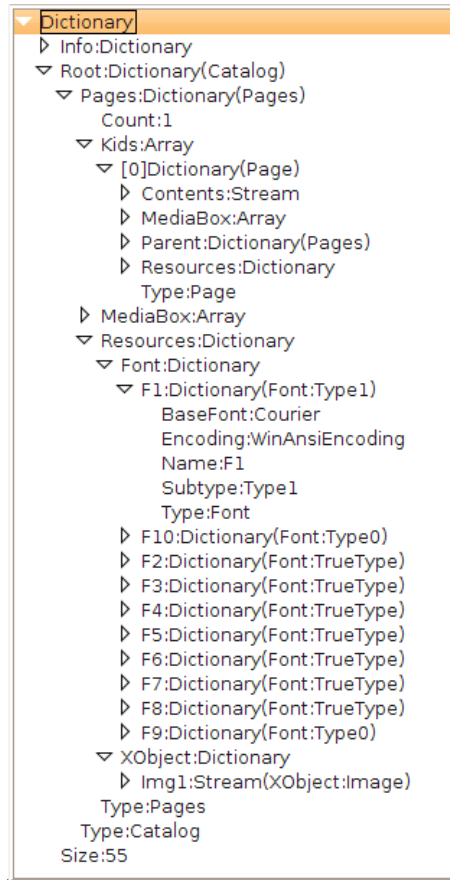


## Portable Document Format (2)

- dictionaries
- streams (command streams or binary data)
- streams are mostly compressed
- encryption support
- consists of several kinds of objects
  - boolean, numeric, string, name, array
  - stream, dictionary



## Portable Document Format (3)



Courier: Toto je pokusný text s češtinou – ěščřžýáíé

Times: Toto je pokusný text s češtinou – ěščřžýáíé

Tahoma: Toto je pokusný text s češtinou – ěščřžýáíé

Arial: Toto je pokusný text s češtinou – ěščřžýáíé

Arial Narrow: Toto je pokusný text s češtinou – ěščřžýáíé

**Arial Black: Toto je pokusný text s češtinou – ěščřžýáíé**

Arial Rounded: Toto je pokusný text s češtinou – ěščřžýáíé

Verdana: Toto je pokusný text s češtinou – ěščřžýáíé

Sans serif: Toto je pokusný text s češtinou – ěščřžýáíé



## Portable Document Format (4)

- low level commands to create content
  - draw text
  - create and draw a path
  - draw an image
- no high level concepts
  - high level text formatting like paragraphs, justification, automatic word-wrapping
  - tables





## Portable Document Format (5)

<ul style="list-style-type: none"><li>▼ Dictionary<ul style="list-style-type: none"><li>▶ ID:Array</li><li>▼ Root:Dictionary(Catalog)<ul style="list-style-type: none"><li>▼ Pages:Dictionary(Pages)<ul style="list-style-type: none"><li>Count:3</li><li>▼ Kids:Array<ul style="list-style-type: none"><li>▶ [0]Dictionary(Page)</li><li>▶ [1]Dictionary(Page)</li><li>▼ [2]Dictionary(Page)<ul style="list-style-type: none"><li>▶ Contents:Stream</li><li>▶ MediaBox:Array</li><li>▶ Parent:Dictionary(Pages)</li><li>▶ Resources:Dictionary<ul style="list-style-type: none"><li>Type:Page</li></ul></li></ul></li><li>Type:Pages</li><li>Type:Catalog</li><li>Version:1.4</li><li>Size:22</li></ul></li></ul></li></ul></li></ul></li></ul>	<pre>/F0 1 Tf BT 12 0 0 12 298 381 Tm ( Hello World - 0 ) Tj 0 18 -18 0 298 381 Tm ( Hello World - 1 ) Tj -24 0 0 -24 298 381 Tm ( Hello World - 2 ) Tj 0 -30 30 0 298 381 Tm ( Hello World - 3 ) Tj ET 1 0 0 RG 119 158 268 357 re s</pre>
--	---

A diagram illustrating the layout of a PDF page. The text is rotated 90 degrees counter-clockwise. The text elements are:

- 2 - Hello World - 0
- 3 - Hello World - 1
- 3 - Hello World - 2
- 3 - Hello World - 3



## Key Features



## Text Extraction

- commandline tool „ExtractText“
- possible options
  - automatic sorting of text chunks
  - different encodings
  - optional selection of start and/or end page
  - simple HTML-output

## Merging/Splitting pdfs

- commandline tool „PDFSplit“
  - splits a pdf at the given page
  - known issue: PDFBox doesn't split the used resources -> results are too large
- commandline tool „PDFMerge“
  - merges two given pdfs into a new one







## Creating Images

- commandline tool „PDFToImage“
  - converts a page to an image
  - supports several output formats
  - supports several color models
  - optional selection of start and/or end page
  - different resolutions
  - supports a different crop box



## Creating a Lucene Index

- commandline tool „IndexFiles“
  - part of the lucene component
  - works on all pdf files in the given directory
  - creates a new index or updates an existing index



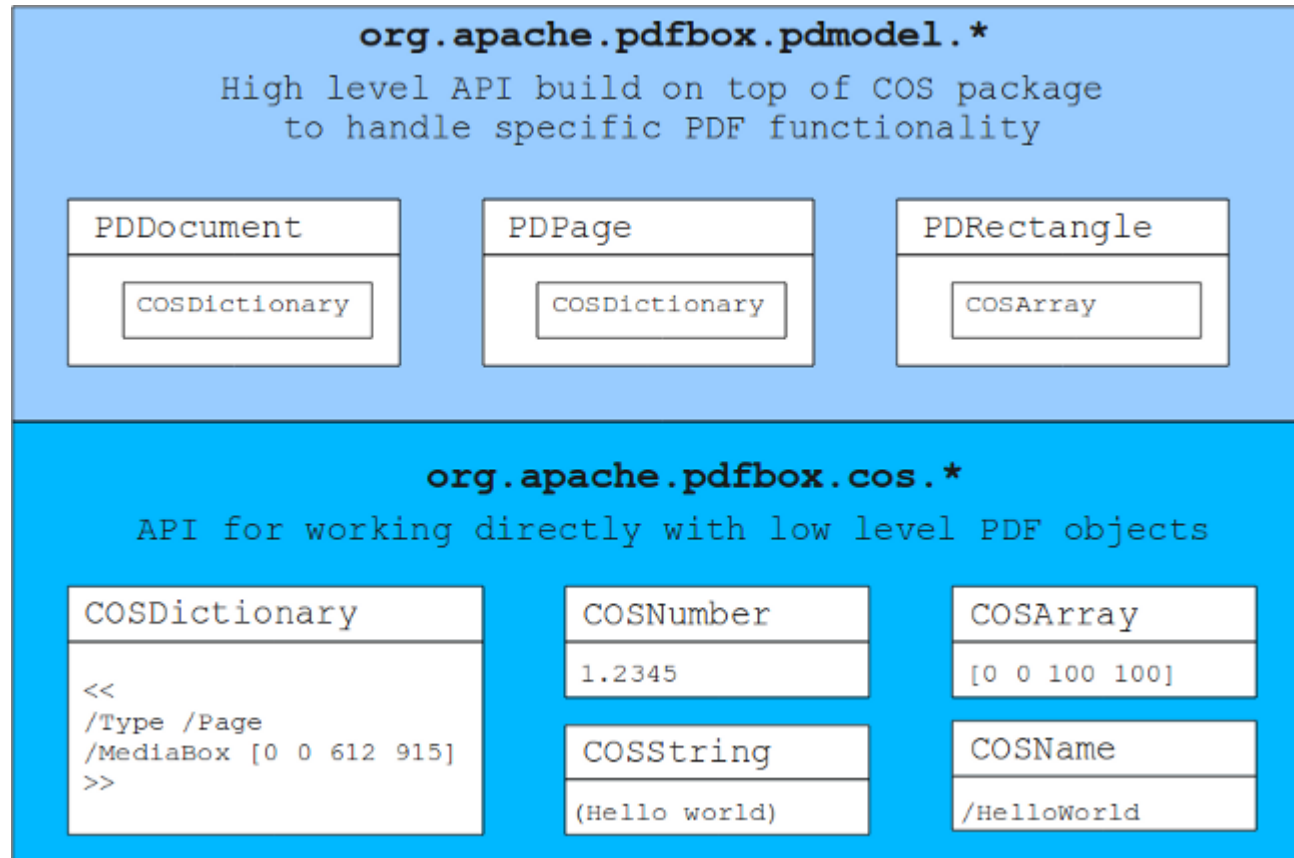
## Printing

- commandline tool „PrintPDF“
  - prints the given pdf
  - printing with/without dialog
  - known limitations: only properly rendered pdfs can be printed



## Using PDFBox Inside Applications

## Memory Presentation Model



## Creating an Empty pdf

```
PDDocument doc = null;
try {
    doc = new PDDocument();
    doc.addPage(new PDPage());
    doc.save("blank.pdf");
}
catch(IOException exception) {
    doc.close();
}
```



## Adding Some Simple Text

```
// create a font
PDFont font = PDType1Font.HELVETICA;
// create a content stream
PDContentStream stream = new PDContentStream(doc,page);
// set font and font size
stream.setFont(font, 12.0f);
// add text
stream.beginText();
stream.moveTextPositionByAmount(100,100);
stream.drawText(„Hello world!!“);
stream.endText();
stream.close();
```



## Using the Text Matrix

```
// create a content stream
PDContentStream stream = new PDContentStream(doc,page);
// set font and font size
stream.setFont(font, 12.0f);
// add text
stream.beginText();
for (int i=0;i<8;i++) {
    stream.setTextRotation(-i*Math.PI*0.25,100,100);
    stream.drawText(", Hello world -" + i);
}
stream.endText();
stream.close();
```







## Using the Text Matrix

Hello World - 0  
Hello World - 1  
Hello World - 2  
Hello World - 3  
Hello World - 4  
Hello World - 5  
Hello World - 6  
Hello World - 7

## Text Related Methods

```
public void setFont(PDFont font, float fontSize)

public void beginText()
public void endText()

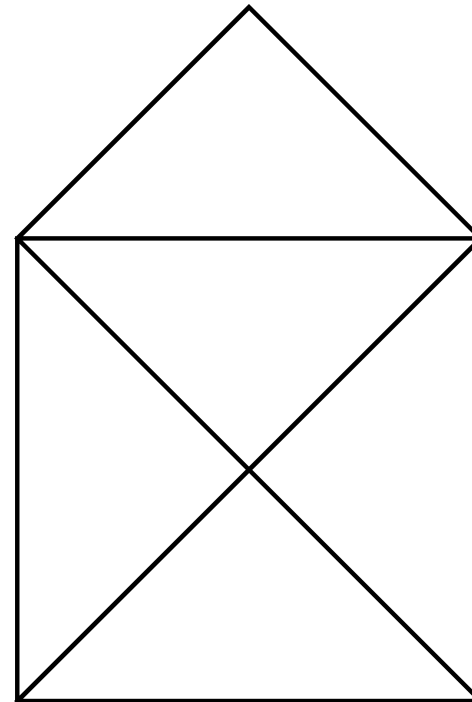
public void moveTextPositionByAmount(float x, float y)

public void setTextMatrix(double a, double b, double c, double d,
    double e, double f);
public void setTextScaling(double sx, double sy,
    double tx, double ty);
public void setTextTranslation(double tx, double ty)
public void setTextRotation(double angle,
    double tx, double ty);
```



## Creating Graphical Content

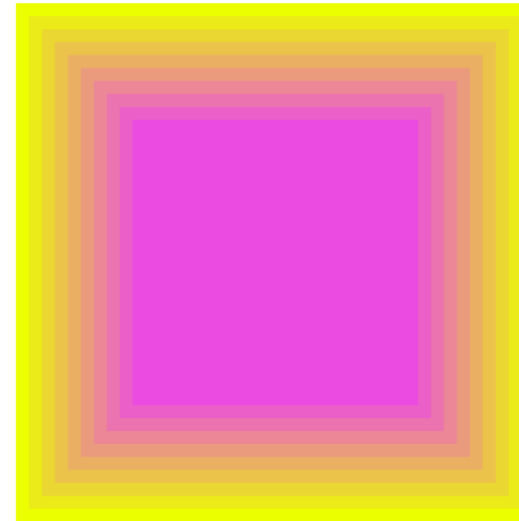
```
PDPageContentStream stream = new PDPageContentStream(doc, page);  
stream.setLineWidth(3);  
stream.moveTo(100, 100);  
stream.lineTo(400, 400);  
stream.lineTo(100, 400);  
stream.lineTo(100, 100);  
stream.lineTo(400, 100);  
stream.lineTo(400, 400);  
stream.lineTo(250, 550);  
stream.lineTo(100, 400);  
stream.lineTo(400, 100);  
stream.stroke();  
stream.close();
```



## Working with Colors

```
PDPageContentStream stream = new PDPageContentStream(doc, page)

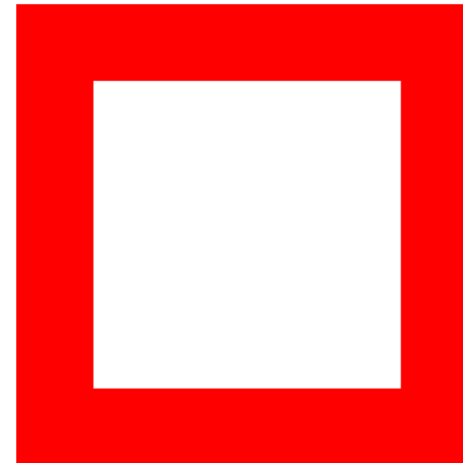
for (int i=0;i<10; i++) {
    stream.setNonStrokingColor(255-20*i, 255-20*i, 25*i);
    stream.fillRect( 10+10*i, 10+10*i, 400-20*i, 400-20*i);
}
stream.close();
```





## Path Filling

```
PDPageContentStream stream = new PDPageContentStream(doc, page);  
stream.setStrokeColor( Color.BLACK);  
stream.setNonStrokingColor( Color.RED);  
stream.addRect(100,100,300,300);  
stream.addRect(150,150,200,200);  
stream.fill(PathIterator.WIND_EVEN_ODD);  
stream.close();
```



## Graphics Related Methods

```
public void setNonStrokingColorSpace(PDColorSpace colorSpace)
public void setStrokingColorSpace( PDColorSpace colorSpace )

public void drawImage(PDXObjectImage image, float x, float y)

public void concatenate2CTM(double a, double b, double c,
    double d, double e, double f)

public void addRect(float x, float y, float width, float height)
public void fillRect(float x, float y, float width, float height)
public void addBezier312(float x1, float y1, float x2, float y2,
    float x3, float y3)
public void addBezier32(float x2, float y2, float x3, float y3)
public void addBezier31(float x1, float y1, float x3, float y3)
public void drawPolygon(float[] x, float[] y)
public void fillPolygon(float[] x, float[] y)
```



## Printing a pdf

```
PDDocument doc = PDDocument.load(„test.pdf“);
PrinterJob job = PrinterJob.getPrinterJob();

// silent print
doc.silentPrint(job);

// open a print dialog
doc.print(jobs)
```



# Converting a pdf into Images

```
// load pdf
PDDocument doc = PDDocument.load(„test.pdf“);
// get all images
List<PDPage> pages = doc.getDocumentCatalog().getAllPages();
// iterate through all pages
Iterator<PDPage> iterator = pages.iterator();
while (iterator.hasNext())
{
    // create a buffered image
    BufferedImage image = iterator.next().convertToImage();
    // save the buffered image, e.g. using ImageIO
    ...
}
```



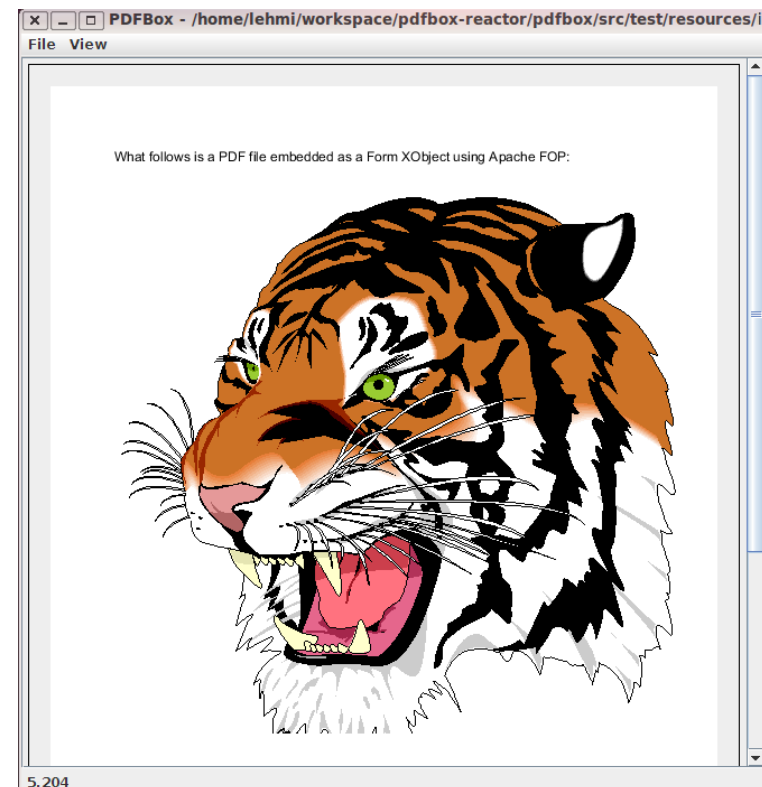
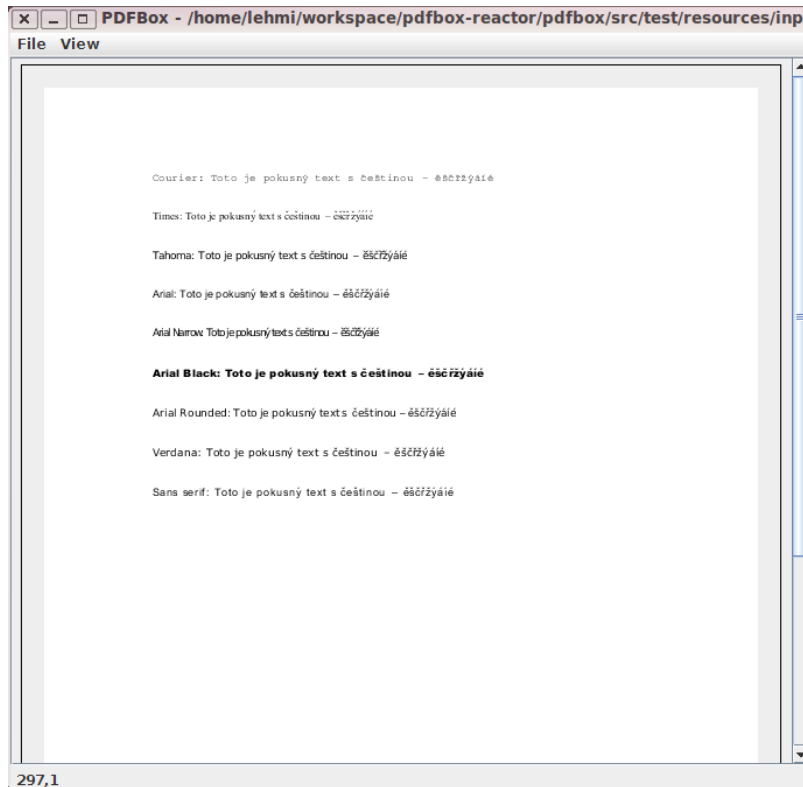




## Debugging pdfs (1)

- PDFReader bundled with PDFBox
- GUI based tool
  - +/- flips pages
  - print
  - no zoom feature (yet)

## Debugging pdfs (2)





## Debugging pdfs (3)

- PDFDebugger bundled with PDFBox
- GUI based tool
- Split-screen
  - left side: treeview of all objects
  - right side: content of the selected object





## Debugging pdfs (5)

- WriteDecodedDoc bundled with PDFBox
- command line tool to decompress pdfs
- decompressed pdf can be easily opened in any editor



## Debugging pdfs (6)

```
sample_fonts_solidconvetor_decompressed.pdf - Kate
Datei Bearbeiten Ansicht Gehe zu Lesezeichen Sitzungen Extras Einstellungen Hilfe
Neu Öffnen Zurück Nach vorne Speichern Speichern unter Schließen Rückgängig Wiederherstellen
Dokumente
Dateisystem-Browser
25.92 -1.405 559.601 791.405 re
W n
Q
Q
q
0 Tr
0.113 w
0 g
q
25.92 -1.405 559.601 791.405 re
W n
99 Tz
BT
/F4 9.3624 Tf
1 0 0 1 92.585 710.476 Tm
[(C-3(o)-3(u)-3(r)-3(i)-3(e)-3(r)-3(:)-3( )-3(T)-3(o)-3(t)-3(o)-3( )-3(j)-3(e)-3( )-3(p)-3(o)-3(k)-3(u)-3(s)-3(n)-3(0)-3( )-3(t
ET
Q
0 g
q
25.92 -1.405 559.601 791.405 re
W n
99 Tz
BT
/F4 9.3624 Tf
1 0 0 1 323.825 710.476 Tm
[(0)]TJ
ET
Q
0 g
q
25.92 -1.405 559.601 791.405 re
W n
99 Tz
Zeile: 1 Spalte: 1 N/L ZEILE sample_fonts_solidconvetor_decompressed.pdf
Terminal In Dateien Suchen
```

## Future Prospects

- Improve font support
  - Type3, CID, encoding ...
- Improve rendering
  - Tiling + shading patterns, colorspaces
- Split to several components
  - Core: parser, text extraction, no gui  
-> GAE, Android...
  - Rendering
  - Examples
- GUI handling all commandline tools





## Get Involved

<http://pdfbox.apache.org>

[users@pdfbox.apache.org](mailto:users@pdfbox.apache.org)  
[dev@pdfbox.apache.org](mailto:dev@pdfbox.apache.org)

<http://www.apache.org/~lehmi/apachecon>

<http://www.apache.org/foundation/getinvolved.html>