working together!

PDF Days Europe 2022 | Berlin

# The PDF Detectives

*The dynamically typed duo is on the case!*

# Who are we?

- Kevin Willems
  - Presales Engineer @ iText


- Michaël Demey
  - Research Manager @iText

# Introducing iText Software

## SDK

### iText 7 Suite

**Best documented and most feature-rich PDF library**

Open-source library for PDF generation and management,

- PDF/A, PDF/UA,
- Digital signing,
- Security,
- HTML conversion,
- Redaction, OCR, etc.,

and closed source add-ons for rendering capabilities, advanced data extraction, advanced forms processing, PDF optimization, advanced language script capabilities, Office to PDF conversion and much more.

## High convenience tools

### iText DITO

**Fast template-based document generation**

High convenience technology block – save time

Based on iText technology

- WYSIWYG Template designer
- Conditional logic for content visibility
- Data binding
- Advanced data visualizations (barcodes/QR codes, charts…)
- PDF/UA Compliance Assistance
- Template Manager with version, dependency and permission control
- Containerized RESTful PDF-producer API (also available as native Java SDK).

# Mission Briefing

# Goal of the Presentation

- Talk about some interesting behavior in PDF files

  - Debug

  - Explain

  - Resolve (hopefully ☺)

# The Toolbelt

# Inventory Report!

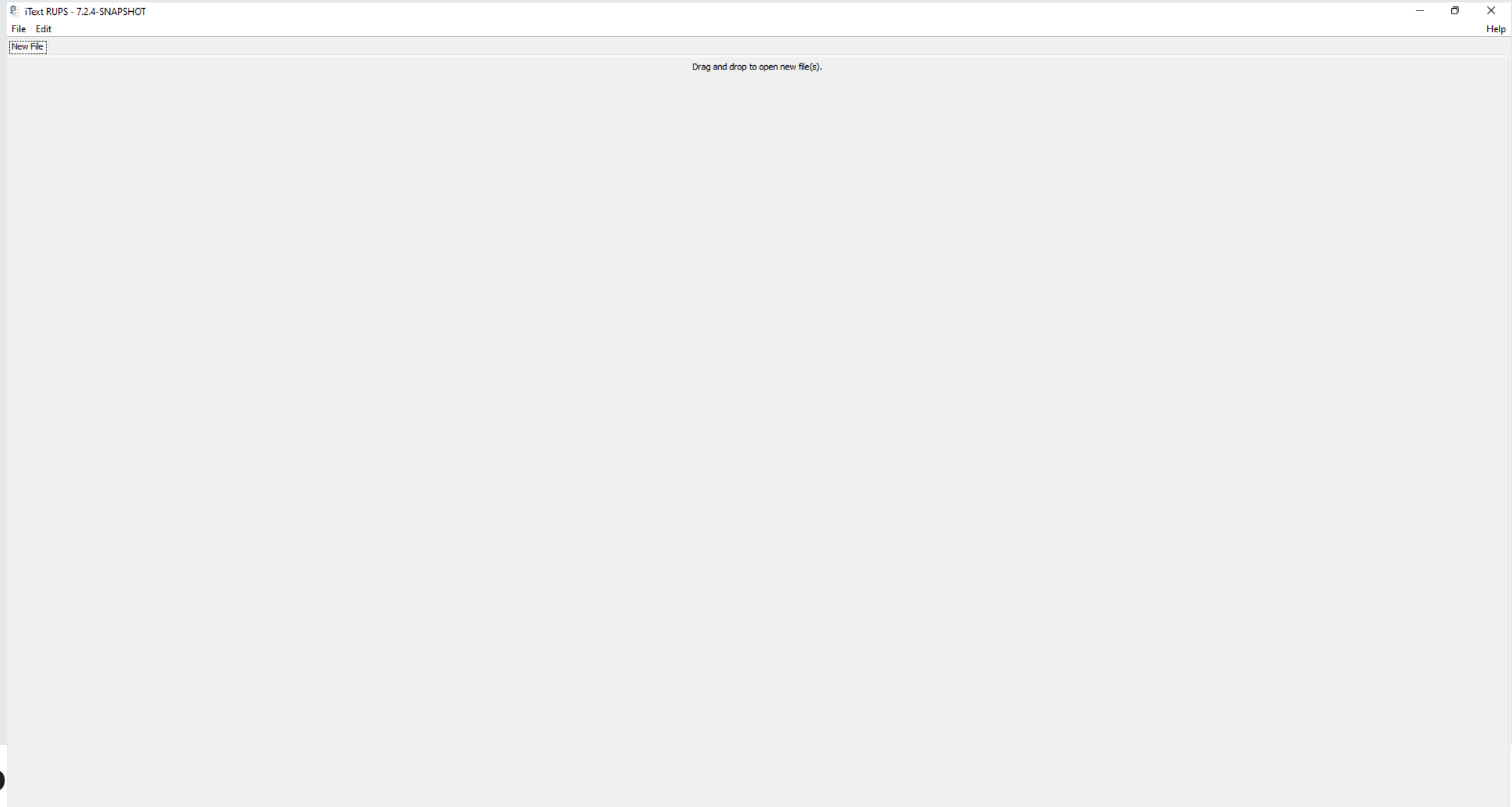Analysis and Debugging Tools

- iText RUPS (+iText Core)
- Acrobat Pro

Accessibility and Tagging

- VeraPDF
- PAC

Low Level Tools

- Hex Editor
- Notepad++

# RUPS "Demo"

File   Edit                                                                                           Help

map-2.pdf

PDF Object Tree (map-2.pdf)
/Info: 1 0 R Modified; -> Dictionary
/Root: 1622 0 R -> Dictionary of type: /Catalog

| Pages | Outlines | Structure | Form | XFA | XREF | Plain Text |

| Object | Page |
| --- | --- |
| Object 2 | Page 1 |
| Object 82 | Page 2 |
| Object 123 | Page 3 |
| Object 130 | Page 4 |
| Object 139 | Page 5 |
| Object 167 | Page 6 |
| Object 198 | Page 7 |
| Object 231 | Page 8 |
| Object 257 | Page 9 |
| Object 268 | Page 10 |
| Object 370 | Page 11 |

Stream   XFA   Debug Info   Console

File   Edit                                                                                                    Help

map-2.pdf

PDF Object Tree (map-2.pdf)
- /Info: 1 0 R Modified; -> Dictionary
- /Root: 1622 0 R -> Dictionary of type: /Catalog
  - Dictionary of type: /Catalog
    - /MarkInfo: Dictionary of type: /MarkInfo
    - /Pages: 397 0 R -> Dictionary of type: /Pages
      - Dictionary of type: /Pages
        - /Count: 11
        - /Kids: Array
          - 395 0 R -> Dictionary of type: /Pages
            - Dictionary of type: /Pages
              - /Count: 8
              - /Kids: Array
                - 2 0 R -> Dictionary of type: /Page
                  - Page 1
                    - /Annots: Array
                    - /Contents: 81 0 R -> Stream
                    - /MediaBox: Array
                    - /Parent: 395 0 R -> Dictionary of type: /Pages
                    - /Resources: Dictionary
                    - /StructParents: 0
                    - /Type: /Page
                - 82 0 R -> Dictionary of type: /Page
                - 123 0 R -> Dictionary of type: /Page
                - 130 0 R -> Dictionary of type: /Page
                - 139 0 R -> Dictionary of type: /Page
                - 167 0 R -> Dictionary of type: /Page
                - 198 0 R -> Dictionary of type: /Page
                - 231 0 R -> Dictionary of type: /Page
              - /Parent: 397 0 R -> Dictionary of type: /Pages
              - /Type: /Pages
          - 396 0 R -> Dictionary of type: /Pages
        - /Type: /Pages
  - /StructTreeRoot: 398 0 R Modified; -> Dictionary of type: /StructTreeRoot
  - /Type: /Catalog

| Pages | Outlines | Structure | Form | XFA | XREF | Plain Text |

| Object | Page |
|---|---|
| Object 2 | Page 1 |
| Object 82 | Page 2 |
| Object 123 | Page 3 |
| Object 130 | Page 4 |
| Object 139 | Page 5 |
| Object 167 | Page 6 |
| Object 198 | Page 7 |
| Object 231 | Page 8 |
| Object 257 | Page 9 |
| Object 268 | Page 10 |
| Object 370 | Page 11 |

| Key | Value | |
|---|---|---|
| /Annots | [ 42 0 R 43 0 R 44 0 R 45 0 R 46 0 R... | ✖ |
| /Contents | 81 0 R | ✖ |
| /MediaBox | [ 0 0 792 612 ] | ✖ |
| /Parent | 395 0 R | ✖ |
| /Resources | << /ExtGState << /G12 12 0 R /G2... | ✖ |
| /StructParents | 0 | ✖ |
| /Type | /Page | ✖ |
| / | | ⊕ |

| Stream | XFA | Debug Info | Console |

File   Edit                                                                                                    Help

map-2.pdf

PDF Object Tree (map-2.pdf)

- /Info: 1 0 R Modified; -> Dictionary
- /Root: 1622 0 R -> Dictionary of type: /Catalog
  - Dictionary of type: /Catalog
    - /MarkInfo: Dictionary of type: /MarkInfo
    - /Pages: 397 0 R -> Dictionary of type: /Pages
      - Dictionary of type: /Pages
        - /Count: 11
        - /Kids: Array
          - 395 0 R -> Dictionary of type: /Pages
            - Dictionary of type: /Pages
              - /Count: 8
              - /Kids: Array
                - 2 0 R -> Dictionary of type: /Page
                  - Page 1
                    - /Annots: Array
                    - /Contents: 81 0 R -> Stream
                      - Stream
                        - /Filter: /FlateDecode
                        - /Length: 9385
                    - /MediaBox: Array
                    - /Parent: 395 0 R -> Dictionary of type: /Pages
                    - /Resources: Dictionary
                    - /StructParents: 0
                    - /Type: /Page
                - 82 0 R -> Dictionary of type: /Page
                - 123 0 R -> Dictionary of type: /Page
                - 130 0 R -> Dictionary of type: /Page
                - 139 0 R -> Dictionary of type: /Page
                - 167 0 R -> Dictionary of type: /Page
                - 198 0 R -> Dictionary of type: /Page
                - 231 0 R -> Dictionary of type: /Page
              - /Parent: 397 0 R -> Dictionary of type: /Pages
              - /Type: /Pages
          - 396 0 R -> Dictionary of type: /Pages
        - /Type: /Pages
  - /StructTreeRoot: 398 0 R Modified; -> Dictionary of type: /StructTreeRoot

Pages   Outlines   Structure   Form   XFA   XREF   Plain Text

| Object | Page |
| --- | --- |
| Object 2 | Page 1 |
| Object 82 | Page 2 |
| Object 123 | Page 3 |
| Object 130 | Page 4 |
| Object 139 | Page 5 |
| Object 167 | Page 6 |
| Object 198 | Page 7 |
| Object 231 | Page 8 |
| Object 257 | Page 9 |
| Object 268 | Page 10 |
| Object 370 | Page 11 |

| Key | Value | |
| --- | --- | --- |
| /Filter | /FlateDecode | ✕ |
| /Length | 9385 | ✕ |
| / | | ⊕ |

Stream   XFA   Debug Info   Console

```
.23999999 0 0 -.23999999 0 612 cm
q
0 0 3301 2551 re
W*
n
q
4.1666665 0 0 4.1666665 0 0 cm
1 1 1 RG
1 1 1 rg
/G3 gs
0 0 1056 816 re
f
Q
Q
q
0 62.499996 362.5 37.500004 re
W*
n
q
4.1666665 0 0 4.1666665 0 0 cm
/G3 gs
BT
/F4 9 Tf
```

PDF association

| Number | Object | Byte Offset |
|---|---|---|
| 1 | Dictionary of type: /Page | 393729 |
| 2 | Indirect object | 393729 |
| 3 | Indirect object | 398265 |
| 4 | Indirect object | 403919 |
| 5 | Indirect object | 406015 |
| 6 | Indirect object | 406360 |
| 7 | Indirect object | 406794 |
| 8 | Indirect object | 408253 |
| 9 | Indirect object | 408545 |
| 10 | Indirect object | 408925 |
| 11 | Indirect object | 410680 |
| 12 | Indirect object | 764116 |
| 13 | Indirect object | 766456 |
| 14 | Indirect object | 789653 |
| 15 | Indirect object | 789701 |
| 16 | Indirect object | 790215 |
| 17 | Dictionary of type: /Outlines | 790620 |
| 18 | Indirect object | 790996 |
| 19 | Indirect object | 791039 |
| 20 | Indirect object | 793419 |
| 21 | Indirect object | 794997 |
| 22 | Indirect object | 796467 |
| 23 | Indirect object | 797809 |
| 24 | Indirect object | 798356 |
| 25 | Dictionary of type: /Pages | 800562 |
| 26 | Stream of type: /Metadata | 804759 |
| 27 | Dictionary | 804819 |
| 28 | Dictionary of type: /StructTreeRoot | 807378 |
| 29 | Indirect object | Obj. Stream #18 ( 791039 ) |
| 30 | Indirect object | Obj. Stream #18 ( 791039 ) |
| 31 | Indirect object | Obj. Stream #18 ( 791039 ) |
| 32 | Indirect object | Obj. Stream #18 ( 791039 ) |
| 33 | Indirect object | Obj. Stream #18 ( 791039 ) |
| 34 | Indirect object | Obj. Stream #18 ( 791039 ) |
| 35 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 36 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 37 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 38 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 39 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 40 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 41 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 42 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 43 | Dictionary | Obj. Stream #18 ( 791039 ) |
| 44 | Dictionary | Obj. Stream #18 ( 791039 ) |

# iText RUPS

- It's on GitHub ☺

# Inventory Report!

Analysis and Debugging Tools

- iText RUPS (+iText Core)
- Acrobat Pro

Accessibility and Tagging

- VeraPDF
- PAC

Low Level Tools

- Hex Editor
- Notepad++

# Cases

# Case 1 - Usecase



⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚

⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚⬚

This text was created with a type 1 font

# Case 1 - Usecase

# Case 1 - Debugging

working together!



PDF association

# Case 1 - Debugging

# Case 1 - Debugging

A **type 0 font** is a virtual font that references multiple component fonts.

# Case 1 - Debugging

A **type 0 font** is a virtual font that references multiple component fonts.

| Key | Type | Value |
|---|---|---|
| | | Table 121 – Entries in a Type 0 font dictionary (continued) |
| DescendantFonts | array | *(Required)* A one-element array specifying the CIDFont dictionary that is the descendant of this Type 0 font. |
| ToUnicode | stream | *(Optional)* A stream containing a CMap file that maps character codes to Unicode values (see 9.10, "Extraction of Text Content"). |

**CID fonts are special fonts designed to display languages with more than 256 characters.**

The main features that CID fonts add are the ability to have 16bit values (so 65535 separate CID characters rather than 256) and much more sophisticated and more flexible unicode settings for extraction. **Predefined CMAPs (or custom ones embedded by the user) allow for text extraction to provide appropriate values.**

A **type0 font** is a virtual font that references multiple component fonts.

Table 121 – Entries in a Type 0 font dictionary (continued)

| Key | Type | Value |
| --- | --- | --- |
| DescendantFonts | array | *(Required)* A one-element array specifying the CIDFont dictionary that is the descendant of this Type 0 font. |
| ToUnicode | stream | *(Optional)* A stream containing a CMap file that maps character codes to Unicode values (see 9.10, "Extraction of Text Content"). |

EXAMPLE 1    This example illustrates a Type 0 font that uses the Identity-H CMap to map from character codes to CIDs and whose descendant CIDFont uses the Identity mapping from CIDs to TrueType glyph indices. Text strings shown using this font simply use a 2-byte glyph index for each glyph. In the absence of a ToUnicode entry, no information would be available about what the glyphs mean.

```
14  0  obj
    <<  /Type  /Font
        /Subtype  /Type0
        /BaseFont  /Ryumin–Light
        /Encoding  /Identity–H
        /DescendantFonts  [15 0 R]
        /ToUnicode  16 0 R
```

**EXAMPLE 1**     This example illustrates a Type 0 font that uses the Identity-H CMap to map from character codes to CIDs and whose descendant CIDFont uses the Identity mapping from CIDs to TrueType glyph indices. Text strings shown using this font simply use a 2-byte glyph index for each glyph. In the absence of a ToUnicode entry, no information would be available about what the glyphs mean.

```
14  0  obj
    <<  /Type  /Font
        /Subtype  /Type0
        /BaseFont  /Ryumin–Light
        /Encoding  /Identity–H
        /DescendantFonts  [15 0 R]
        /ToUnicode  16 0 R
```

/F1: 6 0 R -> Dictionary of type: /Font
 Dictionary of type: /Font
  /BaseFont: /YRCJLH+SourceHanSerif-Regular-Identity-H
  /DescendantFonts: Array
  /Encoding: /Identity-H
  /Subtype: /Type0
  /Type: /Font

# Case 1 - Solution

# Case 2 - Usecase

Case 2 - Debugging

# Case 2 - Debugging



Table 87 – XObject Operator

| Operands | Operator | Description |
|----------|----------|-------------|
| *name* | **Do** | Paint the specified XObject. The operand *name* shall appear as a key in the **XObject** subdictionary of the current resource dictionary (see 7.8.3, "Resource Dictionaries"). The associated value shall be a stream whose **Type** entry, if present, is **XObject**. The effect of **Do** depends on the value of the XObject's **Subtype** entry, which may be **Image** (see 8.9.5, "Image Dictionaries"), **Form** (see 8.10, "Form XObjects"), or **PS** (see 8.8.2, "PostScript XObjects"). |



```
1   %PDF-1.7
2   %âãÏÓ
3   5 0 obj
4   <</Length 475>>stream
5   q
6   0.80553 0 0 0.80553 36 155.15 cm
7   /Fm1 Do
8   Q
9   q
10  BT
11  /F1 13 Tf
12  57.1 136 Td
13  <000401760003017d016f011a00030190015d016f011e0176019a00030189017d0176011a>Tj
14  ET
15  Q
16  q
17  BT
18  /F1 13 Tf
19  49.22 106.59 Td
20  <000400030128018c017d01500003016901b50175018901900003015d0176019a017d0003019a015a011e>Tj
21  <0003>Tj
22  ET
23  Q
24  q
25  BT
26  /F1 13 Tf
27  85.45 85.18 Td
28  <0189017d0176011a0376>Tj
29  ET
30  Q
31  q
32  BT
33  /F1 13 Tf
34  48.6 55.77 Td
35  <005e0189016f01020190015a034a0003005e015d016f011e01760110011e0003010201500102015d01760358>Tj
36  ET
37  Q
```

# Case 2 - Debugging

# Case 2 - Debugging

# Case 2 - Debugging

# Case 2 - Solution

# Case 3 - Usecase

# Case 3 - Debugging

# Case 3 – Solution?

# Case 3 – Solution!

working together!



PDF association

# Case 4

Tim Allison 🚀
@_tallison

Calling PDF devs, both commercial and open source, take a look at the file attached to:

issues.apache.org/jira/browse/PD...

This file causes infinite loops in a bunch of tools that I thought would handle this more robustly.
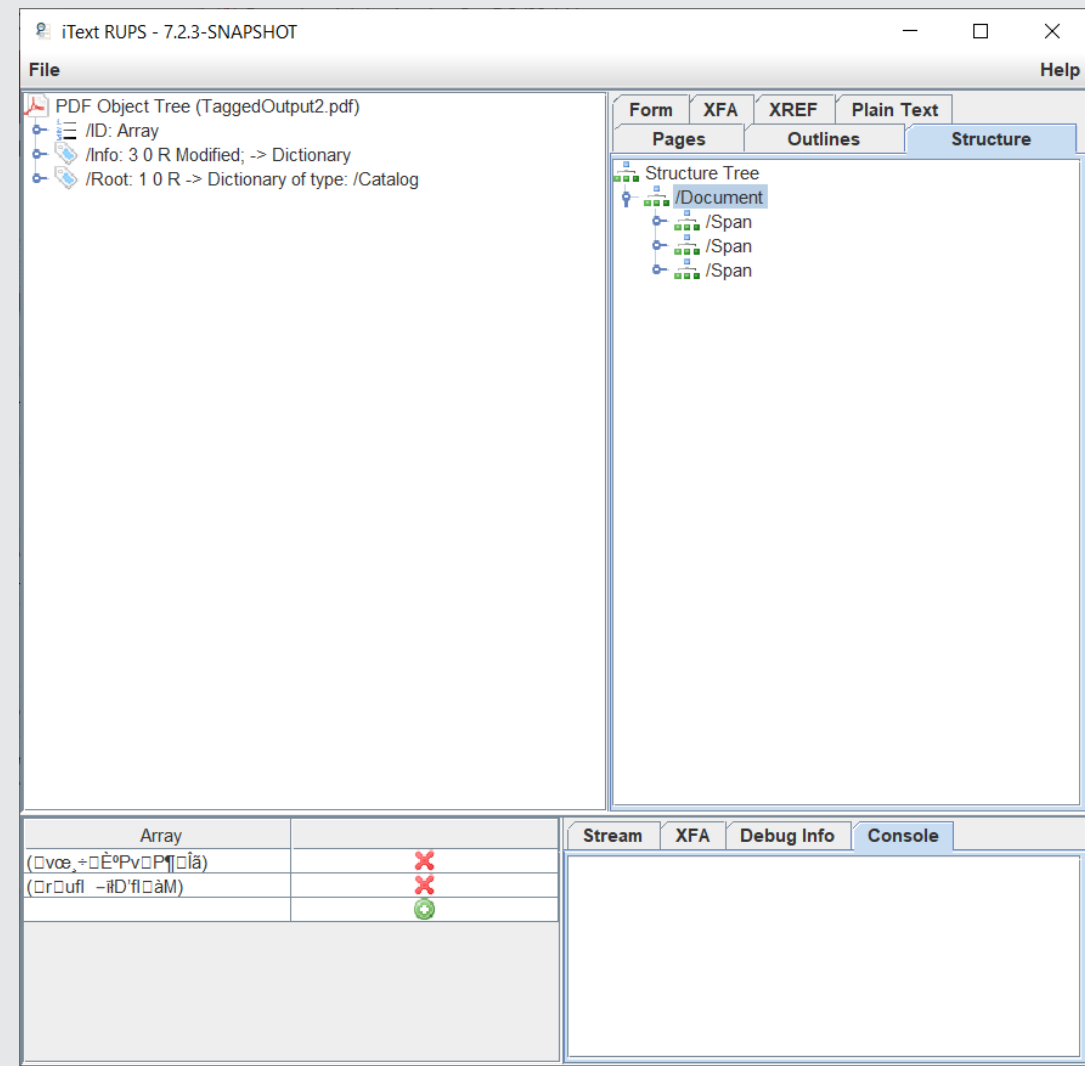
8:29 PM · Apr 12, 2022 · Twitter Web App

PDFBox / PDFBOX-5415

## Infinite loop in ExtractText in 2.x branch on a specific pdf

**Description**

DavidAvant reported an infinite loop in Tika and provided an example file. I can reproduce this with the latest PDFBox app 2.0.26-SNAPSHOT's ExtractText.

File: https://issues.apache.org/jira/secure/attachment/13042292/map.pdf

Adobe and a slightly out of date pdftotext also have problems with this file.

# Case 4 – The initial look

# Case 4 – The initial look



It crashes the viewer

# Case 4- Break out the tools

- Reproduce the issue as-is
  - Run PDF Box code
  - Encounter reported behavior

# Case 4 – Switch up the tools

- Run the same code on iText 7
  - Is the bug only in PDF Box or is it a file issue
  - Encounter the same behavior

# Case 4 – Dissecting the body

```
.23999999 0 0 -.23999999 0 612 cm
/G3 gs
/G336 gs
/X334 Do
/G8 gs
/G337 gs
/X334 Do
/G8 gs
```

```
.23999999 0 0 -.23999999 0 612 cm
/G3 gs
/G332 gs
/X330 Do
/G8 gs
/G333 gs
/X330 Do
/G8 gs
```

```
.23999999 0 0 -.23999999 0 612 cm
/G3 gs
/G328 gs
/X326 Do
/G8 gs
/G329 gs
/X326 Do
/G8 gs
```

working together!

```
    Jenkinsfile                         19  ▶     public static void main(String[] args) throws Exception {
    LICENSE.md                          20             PdfDocument pdfDoc = new PdfDocument(new PdfReader( filename: "page10.pdf"));
    page10.pdf                          21             System.out.println(countXobjDepth(pdfDoc.getFirstPage().getResources()));
m   pom.xml                             22         }
    README.md                           23         public static int countXobjDepth(PdfResources res) {
    root.iml                            24             if(res == null || !res.getPdfObject().containsKey(PdfName.XObject)) {
    ruby-test.pdf                       25                 return 0;
    sonar-project.properties            26             }
    testout.pdf                         27             int curMax = 0;
    Vagrantfile                         28             for(PdfObject child : res.getPdfObject().getAsDictionary(PdfName.XObject).values()) {
    weirdxmp.pdf                        29                 PdfStream stm = (PdfStream) child;
  External Libraries                    30                 int childDepth;
  Scratches and Consoles               31                 if(PdfName.Form.equals(stm.get(PdfName.Subtype))) {
∨   Extensions                          32                     childDepth = 1 + countXobjDepth(new PdfFormXObject(stm).getResources());
  ∨   Java                              33                 } else {
          predefinedExternalAnnotations.j  34                     childDepth = 1; // image -> depth 1
          predefinedExternalAnnotations.j
```

```
Scratch ×
"C:\Program Files\Java\jdk1.8.0_251\bin\java.exe" ...
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#StaticLoggerBinder for further details.
24

Process finished with exit code 0
```

PDF association

50

# Case 4 – Conclusion

- The PDF file isn't broken
- Implementations aren't broken

- PDF files can have complex and intricate constructs
  - Flexibility of the spec

- Parsers need to find ways around this
  - Cache XObjects
  - Lazy Loading of XObjects

# Case 5 – The One with Crowd-sourcing

- Customer reports "broken text extraction"
  - The text was perfectly legible in viewers

# Case 5

- Customer reports "broken text extraction"
  - The text was perfectly legible in viewers
  - When extracting, only some letters would be incorrect
    - Yet, consistent incorrect

# Case 5

- Customer reports "broken text extraction"
  - The text was perfectly legible in viewers
  - When extracting, only some letters would be incorrect
    - Yet, consistent incorrect
  - Across tools
    - So, not an iText issue

- Customer reports "broken text extraction"
  - The text was perfectly legible in viewers
  - When extracting, only some letters would be incorrect
    - Yet, consistent incorrect
  - Across tools
    - So, not an iText issue
  - Through the same producer

# Case 5 – Resolution

- The mapping was completely off

- Erroneous behavior of that producer?
- Obfuscation?

# Mission Report