working together!

PDF Days Europe 2022 | Berlin

# PDF Optimization *Horror* Stories

Grizzly PDF autopsies

# Datalogics

working together!

- **The Adobe PDF Library** – built on Adobe source code
- **Live support** from PDF technology experts
- **SDKs and command-line tools** for building large applications – great for OEMs, system integrators and enterprise developers
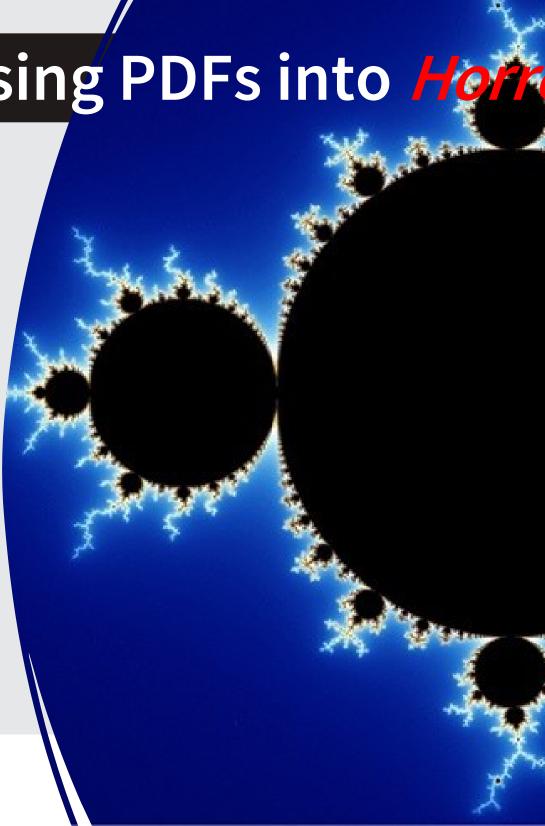
www.datalogics.com

# What makes Processing PDFs into *Horror* stories?

- Ballooning file size
- Hangs /infinite-loops

- Most issues are edge-cases rather than Horror Stories

**PDF association**

# Sub-optimal StructureTrees

- Case 1: The structure tree that grew when it was copied.

- Case 2: Phantom Limbs in the StructureTree

# 36876: the StructureTree that grew…

\\penguin\support\csts\2725\36876\mergepdf2.pdf file size:3923477

| | | |
|---|---:|---:|
| Images: | 1,440,040 | 36.70% |
| Content Streams: | 163,528 | 4.17% |
| XObject Forms: | 82,780 | 2.11% |
| Fonts: | 152,937 | 3.90% |
| Structure Info*: | 1,751,785 | 44.65% |
| Document Overhead*: | 60,467 | 1.54% |
| Cross Reference Table: | 271,940 | 6.93% |
| Total file Size: | 3,923,477 | 100.00% |

\penguin\support\csts\2725\36876\mergepdf_out1564pzs.pdf file size:

| | | |
|---|---:|---:|
| Images: | 1,549,447 | 1.99% |
| Content Streams: | 163,875 | 0.21% |
| XObject Forms: | 30,489 | 0.04% |
| Fonts: | 152,838 | 0.20% |
| Structure Info*: | 54,688,652 | 70.22% |
| Document Overhead*: | 5,403,782 | 6.94% |
| Cross Reference Table: | 15,895,180 | 20.41% |
| Total file Size: | 77,884,263 | 100.00% |

# 40357: Phantom Limbs in StructureTree

support\csts\1000\40357\PDFUA-Reference-04_(Danish_Blind_Association).pdf

| | | |
|---|---:|---:|
| Images: | 2,311,835 | 75.61% |
| Content Streams: | 181,626 | 5.94% |
| XObject Forms: | 32,332 | 1.06% |
| Fonts: | 276,755 | 9.05% |
| Color Spaces*: | 5,389 | 0.18% |
| Extended Graphic States: | 3,138 | 0.10% |
| Shading Info*: | 8,380 | 0.27% |
| Link Annotations: | 5,608 | 0.18% |
| Acrobat Forms*: | 123 | 0.00% |
| Structure Info*: | 90,757 | 2.97% |
| Bookmarks: | 1,674 | 0.05% |
| Named Destinations*: | 779 | 0.03% |
| Document Overhead*: | 195,810 | 6.40% |
| Cross Reference Table: | 2,007 | 0.07% |
| Total file Size: | 3,057,561 | 100.00% |

| | | |
|---|---:|---:|
| Images: | 5,460,194 | 77.00% |
| Content Streams: | 211,051 | 2.98% |
| XObject Forms: | 33,366 | 0.47% |
| Fonts: | 534,428 | 7.54% |
| Color Spaces*: | 20,558 | 0.29% |
| Extended Graphic States: | 29,244 | 0.41% |
| Shading Info*: | 6,378 | 0.09% |
| Link Annotations: | 11,278 | 0.16% |
| Structure Info*: | 466,734 | 6.58% |
| Bookmarks: | 1,648 | 0.02% |
| Named Destinations*: | 797 | 0.01% |
| Document Overhead*: | 270,228 | 3.81% |
| Cross Reference Table: | 45,340 | 0.64% |
| Total file Size: | 7,091,244 | 100.00% |

# 44874: Cyclic resource trees.

- Linearization/Fast Web Preview
  - How it works
  - How it went wrong (for this file)

# Examining the Resource Tree.

# Resource Tree Wheel

- Each Form XObject contained an identical **copy** of a previously shared common Resource tree.

- Gordian knot  cut by…examining which resources each Form XObject actually used.

# 42856: Font Optimization



```
f:\support\csts\3577\42856\wf33558459b0863-8001-obfoscated-fott.pdf
                      Images:      3,535,012      0.26%
             Content Streams:     12,768,259      0.94%
                XObject Forms:    579,792,125     42.49%
                        Fonts:    787,818,561     57.73%
       Extended Graphic States:          255      0.00%
          Cross Reference Table:    2,192,359      0.16%
              Total file Size:  1,364,662,723    100.00%
```

# Font Stats

- **SELECT COUNT(*) from Font where docid=1**

- (416586,)

- **SELECT COUNT(*),subtype from Font where docid=1 and subtype!="Type0" group by subtype**

- (1, 'CIDFontType0')

- (208291, 'CIDFontType2')

- (2, 'Type1')

- **SELECT DISTINCT BASEFONT from Font where docid=1**

- ('Gotham-Bold',)

- ('NimbusSanL-Regu',)

- ('Gotham-Medium',)

- ('Gotham-Book',)

- ('AllAndNone',)

- ('Helvetica',)

- ('ArialMT',)

- ('Gotham-Light',)

- ('DejaVuSans',)

- **SELECT DISTINCT substr(FONTNAME,7) from Fontdescriptor where FONTNAME like "%+%" and docid=1**

- ('+Gotham-Bold',)

- ('+NimbusSanL-Regu',)

- ('+Gotham-Medium',)

- ('+Gotham-Book',)

- ('+Gotham-Light',)

- ('+DejaVuSans',)

- **select count(OBJECTID),basefont from FONT where docid=1 and subtype!="Type0" group by basefont**

- (1, 'AllAndNone')

- (1, 'ArialMT')

- (3244, 'DejaVuSans')

- (50230, 'Gotham-Bold')

- (50586, 'Gotham-Book')

- (3771, 'Gotham-Light')

- (50230, 'Gotham-Medium')

- (1, 'Helvetica')
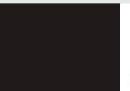
- (50230, 'NimbusSanL-Regu')

# Take-aways

- PDFs are not (always) Trees.

  - Tree nodes aren't referenced more than once

- PDFs are databases

  - of indirect objects

  - of serialized data structures

- Standard PDF processing can sometimes lead to Horrible Files.

  - But not that often.

# Thank you!