



PDF Days Europe 2025

PDF XMP Metadata Validation with Relax NG Schemas and the PoDoFo Library

An Open Source Implementation of ISO 16684-2:2014

<u>Francesco Pretto</u> | Software Engineer | The PoDoFo Community

Outline

- Background/context
 - XMP, PoDoFo library, problem definition
- Normative framework
 - Requirements and interesting findings
- The solution development process
- Next Steps/TODOs
- Summary
- Questions



Context - XMP Extensible Metadata Platform

- An ISO standardized cross format (PDF, JPEG, etc.) container for storing metadata
- Developed by Adobe in 2001, initially with focus on PDF
- XML based: still good (focus in formal validation) but old



```
<?xpacket begin="*\phi" id="W5M0MpCehiHzreSzNTczkc9d"?>
    <x:xmpmeta xmlns:x="adobe:ns:meta/" >
      <rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">
        <rdf:Description xmlns:dc="http://purl.org/dc/elements/1.1/"
                         xmlns:xmp="http://ns.adobe.com/xap/1.0/" rdf:about="">
          <dc:title>
            <rdf:Alt>
              <rdf:li xml:lang="x-default">Document title</rdf:li>
            </rdf:Alt>
10
          </dc:title>
11
          <xmp:CreatorTool>Adobe Illustrator CS3</xmp:CreatorTool>
          <xmp:CreateDate>2015-05-21T13:22:16+02:00/xmp:CreateDate>
12
13
          <xmp:ModifyDate>2015-06-12T18:02:02+02:00</xmp:ModifyDate>
        </rdf:Description>
14
15
      </rdf:RDF>
    </x:xmpmeta>
    <?xpacket end="w"?>
```



Context - Issues with XMP and normal validation strategies

XMP has an ambiguous data model (same content, multiple representations):

```
<xmpMM:History>
    <xmpMM:History>
                                                                                         <rdf:Seq>
       <rdf:Seq>
                                                                                           <rdf:li>
          <rdf:li rdf:parseType="Resource">
                                                                                             <rdf:Description>
             <stEvt:action>converted</stEvt:action>
                                                                                               <stEvt:action>converted</stEvt:action>
             <stEvt:instanceID>uuid:f6227579...</stEvt:instanceID>
                                                                                               <stEvt:instanceID>uuid:f6227579...</stEvt:instanceID>
             <stEvt:parameters>converted to PDF/A-1b</stEvt:parameters>
                                                                                               <stEvt:parameters>converted to PDF/A-1b</stEvt:parameters>
             <stEvt:softwareAgent>Preflight</stEvt:softwareAgent>
                                                                                               <stEvt:softwareAgent>Preflight</stEvt:softwareAgent>
             <stEvt:when>2022-06-20T10:56:38+02:00</stEvt:when>
                                                                                               <stEvt:when>2022-06-20T10:56:38+02:00</stEvt:when>
                                                                                   10
                                                                                             </rdf:Description>
          </rdf:li>
                                                                                           </rdf:li>
10
       </rdf:Seg>
                                                                                         </rdf:Seq>
   </xmpMM:History>
                                                                                      </xmpMM:History>
```

- Regular schema based validation is not normally possible
- Adobe <u>XMP Toolkit SDK</u> comes to rescue and is available for C++ (and Java, but it's not open source anymore)
- Validation with XMP data model parsing libraries is possible (<u>veraPDF</u> does it) but arguably less formal than a schema based validation
- Actually I didn't know about XMP Toolkit SDK until 2025 :-)



Context - What is PoDoFo?

- A modern C++17 PDF manipulation library
- Multiplatform (Win/Lin/Mac, also working in Android/iOS)
- Open Source, LGPL2+ (planned: MPL2)
- Advanced font/text cleanup and sanitization capabilities, powered by Adobe Font Development Kit for OpenType (AFDKO)
- Originally written by German developer Dominik Seichter
- Currently developed/maintained by the author of this presentation
- If you are interested, please reach out at:



https://github.com/podofo/podofo



Context - Original problem and other use cases

- Use case: PDF/A reduction, aka "conversion" (with PoDoFo). Existing metadata packets must be sanitized from invalid properties
- Rigorous XMP packet validation (can we do a better job than XMP data model parsing strategy?)
- XMP extensions pluggable validation: <pdfaExtension> in PDF/A-1 to 3, AFRelationship in PDF/A-4
- Platform/language independent validation: RELAX NG validators are available in multiple platforms (at least Java, C, C#/.NET)

```
<?xpacket begin="*\phi" id="W5M0MpCehiHzreSzNTczkc9d"?>
          <x:xmpmeta xmlns:x="adobe:ns:meta/">
            <rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">
              <rdf:Description xmlns:dc="http://purl.org/dc/elements/1.1/"</pre>
               xmlns:xmp="http://ns.adobe.com/xap/1.0/"
               xmlns:pdf="http://ns.adobe.com/pdf/1.3/" rdf:about="">
               <dc:format>application/pdf</dc:format>
               <xmp:CreatorTool>Apache FOP</xmp:CreatorTool>
               <xmp:CreateDate>2019-03-23T16:51:05Z</xmp:CreateDate>
               <xmp:ModifyDate>2022-05-20T17:19:11+02:00
                <xmp:MetadataDate>2019-03-23T16:51:05Z</xmp:MetadataDate>
             <pd:Trapped>True</pdf:Trapped></pd>
      13
               <CustomProp>Custom property value</CustomProp>
              </rdf:Description>
            </rdf:RDF>
          </x:xmpmeta>
          <?xpacket end="r"?>
                                             Requires
Not available
                                             <pd><pdfaExtension>
until PDF2.0
                                             in PDF/A-1 to 3
(PDF/A-4)
```



Normative framework - The requirements

- ISO 16684-1:2019:
 - Graphic technology **Extensible metadata platform (XMP)**, Part 1: Data model, serialization and core properties
- XMP Specification Part 2 [various revisions]: Additional Properties
- ISO 19005 (aka <u>PDF/A</u>):
 Document management Electronic document file format for long-term preservation
 - <u>ISO 19005-1:2005</u> Part 1: **Use of PDF 1.4** (PDF/A-1)
 - <u>ISO 19005-2:2011</u> Part 2: **Use of ISO 32000-1** (PDF/A-2)

 <u>ISO 19005-3:2012</u> Part 3: **Use of ISO 32000-1 with support for embedded files** (PDF/A-3)
 - <u>ISO 19005-4:2020</u> Part 4: **Use of ISO 32000-2** (PDF/A-4)



Normative framework - Standards to the rescue

- ISO 16684-2:2014:
 Graphic technology Extensible metadata platform (XMP), Part 2:
 Description of XMP schemas using RELAX NG
- ISO/IEC 19757-2:2008:
 Information technology Document Schema Definition
 Language (DSDL)
 Part 2: Regular-grammar-based
 validation RELAX NG

 RELAX NG: simple, intuitive and precise validation schema syntax for XML content



ISO 16884-2 sketches the solution for XMP validation but "leaves it as a trivial exercise to the reader"



Solution - Fully implement ISO 16884-2

- 1. Normalize the XMP packets with a canonization algorithm (chapter 5 "Canonical serialization of XMP" of ISO 16684-2:2014)
- 2. Write a full schema for the PDF XMP packets. PDF/A compliances are more accurate with regard to XMP specification revision: https://github.com/ceztko/XMP-RNG-Schema
- 3. [C/C++] Implement a better streaming validation in libXML2: https://gitlab.gnome.org/GNOME/libxml2/-/merge_requests/334 [Merged!] Added xmlRelaxNGValidCtxtClearErrors to recover from errors during streaming validation



Solution - XMP canonization algorithm 1

- Implemented in <u>PoDoFo</u> (C++)
- <rdf:RDF> outermost XML element, which shall contain a single <rdf:Description> element to contain all XMP properties
- Property (and structure fields) attribute notation shall not be used:

```
<rdf:Description rdf:about="" xmp:Rating="3" />
```





Solution - XMP canonization algorithm 2

The rdf:parseType="resource" attribute notation shall not be used





The canonical serialization of XMP shall use a nested rdf:Bag, rdf:Seq, or rdf:Alt element for array values [unclear]:

```
maybe malformed (x) <dc:date>2019-03-23T16:51:05Z</dc:date>?
```



Solution - XMP canonization algorithm 3

 All general qualifiers shall be serialized as XML elements within an rdf:Description element [...]. The value of that XMP entity shall be within an rdf:value element within that rdf:Description element.





Solution - PDF XMP packets schema design

- The schema must be different depending on XMP specification revision (or PDF/A compliance)
- XSLT compatible expressions are introduced to allow for easy templatization of the schema
- Some edge cases present while reproducing XMP complexity with RELAX NG schemas (eg. User defined structures, open choice of lists, XPath validation)

```
1 <rng:element name="rdf:RDF">
         <rng:element name="rdf:Description">
             <rng:attribute name="rdf:about"/>
             <rng:interleave>
                 <rng:ref name="XMP_Properties-dc"/>
                 <rng:ref name="XMP_Properties-xmp"/>
                 <rng:ref name="XMP Properties-xmpMM"/>
                 <!-- ... -->
                 <rng:ref name="XMP_Properties-pdf"/>
10
                 <rng:ref name="XMP Properties-photoshop"/>
11
                 <rng:ref name="XMP Properties-tiff"/>
12
                 <rng:ref name="XMP_Properties-crs"</pre>
13
                     condition="$IsPDFA2OrGreater"/>
14
                 <rng:ref name="XMP Properties-exif"/>
15
                 <rng:ref name="XMP Properties-exif aux"</pre>
                     condition="$IsPDFA2 or $IsPDFA3"/>
16
                 <rng:ref name="XMP Properties-exifEX"</pre>
17
18
                     condition="$IsPDFA40rGreater"/>
19
                 <!-- ... -->
                 <rng:ref name="XMP Properties-pdfaid"/>
20
21
                 <rng:ref name="XMP Properties-pdfaExtension"</pre>
22
                     condition="$IsPDFA1 or $IsPDFA2 or $IsPDFA3"/>
23
             </rng:interleave>
         </rng:element>
24
    </rng:element>
```



Solution - API Integration (C++, PoDoFo)

```
auto packet = PdfXMPPacket::Create(buffer);
packet->PruneAndValidate(PdfALevel::L2B, [](const PdfXMPProperty& prop) {
    cout << "Invalid property \"" << prop.GetName() << "\"" << endl;
});
cout << packet->ToString() << endl;</pre>
```

```
<?xpacket begin="**\phi" id="W5M0MpCehiHzreSzNTczkc9d"?>
    <x:xmpmeta xmlns:x="adobe:ns:meta/">
      <rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">
        <rdf:Description xmlns:dc="http://purl.org/dc/elements/1.1/"
          xmlns:xmp="http://ns.adobe.com/xap/1.0/"
          xmlns:pdf="http://ns.adobe.com/pdf/1.3/" rdf:about="">
          <dc:format>application/pdf</dc:format>
          <xmp:CreatorTool>Apache FOP</xmp:CreatorTool>
9
          <xmp:CreateDate>2019-03-23T16:51:05Z</xmp:CreateDate>
10
          <xmp:ModifvDate>2022-05-20T17:19:11+02:00/xmp:ModifvDate>
11
          <xmp:MetadataDate>2019-03-23T16:51:05Z</xmp:MetadataDate>
12
          <pdf:Trapped>True</pdf:Trapped>
          <CustomProp>Custom property value</CustomProp>
13
14
        </rdf:Description>
15
      </rdf:RDF>
    </x:xmpmeta>
   <?xpacket end="r"?>
```

```
<?xpacket begin="**O" id="W5M0MpCehiHzreSzNTczkc9d"?>
    <x:xmpmeta xmlns:x="adobe:ns:meta/">
      <rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">
        <rdf:Description xmlns:dc="http://purl.org/dc/elements/1.1/"</pre>
          xmlns:xmp="http://ns.adobe.com/xap/1.0/"
          xmlns:pdf="http://ns.adobe.com/pdf/1.3/" rdf:about="">
          <dc:format>application/pdf</dc:format>
          <xmp:CreatorTool>Apache FOP</xmp:CreatorTool>
          <xmp:CreateDate>2019-03-23T16:51:05Z</xmp:CreateDate>
10
          <xmp:ModifvDate>2022-05-20T17:19:11+02:00</xmp:ModifvDate>
11
          <xmp:MetadataDate>2019-03-23T16:51:05Z</xmp:MetadataDate>
12
        </rdf:Description>
13
      </rdf:RDF>
    </x:xmpmeta>
   <?xpacket end="w"?>
```



Bonus - Validate packet with both PDF/A and PDF/UA compliance

■ PDF/UA, PDF/VT, PDF/X properties require extensions in PDF/A-1 to 3

PoDoFo automatically add the extensions

```
<?xpacket begin="�" id="W5M0MpCehiHzreSzNTczkc9d"?>
    <x:xmpmeta xmlns:x="adobe:ns:meta/">
      <rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">
        <rdf:Description xmlns:dc="http://purl.org/dc/elements/1.1/"</pre>
            xmlns:xmp="http://ns.adobe.com/xap/1.0/"
            xmlns:pdf="http://ns.adobe.com/pdf/1.3/"
            xmlns:pdfaid="http://www.aiim.org/pdfa/ns/id/"
            xmlns:pdfuaid="http://www.aiim.org/pdfua/ns/id/" rdf:about="">
          <dc:format>application/pdf</dc:format>
          <xmp:CreatorTool>Apache FOP</xmp:CreatorTool>
11
          <xmp:CreateDate>2019-03-23T16:51:05Z</xmp:CreateDate>
          <xmp:ModifyDate>2022-05-20T17:19:11+02:00</xmp:ModifyDate>
12
13
          <xmp:MetadataDate>2019-03-23T16:51:05Z</xmp:MetadataDate>
          <pdfaid:part>2</pdfaid:part>
14
          <pdfaid:conformance>B</pdfaid:conformance>
15
          <pdfuaid:part>1</pdfuaid:part>
16
17
        </rdf:Description>
18
      </rdf:RDF>
19
    </x:xmpmeta>
    <?xpacket end="r"?>
```

```
<?xpacket begin="�" id="W5M0MpCehiHzreSzNTczkc9d"?>
    <x:xmpmeta xmlns:x="adobe:ns:meta/">
      <rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#">
        <rdf:Description xmlns:dc="http://purl.org/dc/elements/1.1/" ...
            xmlns:pdfaExtension="http://www.aiim.org/pdfa/ns/extension/"
            xmlns:pdfaSchema="http://www.aiim.org/pdfa/ns/schema#" ... >
          <!-- ... -->
          <pdfaid:part>2</pdfaid:part>
          <pdfaid:conformance>B</pdfaid:conformance>
          <pdfuaid:part>1</pdfuaid:part>
10
          <pdfaExtension:schemas>
11
            <rdf:Bag>
12
13
              <rdf:li>
                <rdf:Description>
14
                  <pdfaSchema:namespaceURI>http://www.aiim.org/pdfua/ns/id/
    </pdfaSchema:namespaceURI>
16
                  <pdfaSchema:prefix>pdfuaid</pdfaSchema:prefix>
                  <pdfaSchema:schema>PDF/UA ID Schema</pdfaSchema:schema>
18
19
                  <!-- ... -->
                </rdf:Description>
20
              </rdf:li>
22
            </rdf:Bag>
23
          </pdfaExtension:schemas>
24
        </rdf:Description>
      </rdf:RDF>
    </x:xmpmeta>
    <?xpacket end="r"?>
```

association

Next Steps

- Support extensions: convert from <pdfaExtension:schemas> to RELAX NG schemas (PDF/A-1 to 3)
- Validate AFRelationship (PDF/A-4)
- Handle edge cases: user defined structures (eg. xmpMM:Pantry), open choice of lists (eg. tiff:YCbCrSubSampling), XPath validation
- Feedback welcome in the XMP RNG Schema (and PoDoFo) project: https://github.com/ceztko/XMP-RNG-Schema



Conclusions

- **Future-proof:** XMP validation with RELAX NG is now a practical alternative
- **Open and free:** Schemas released under a liberal license
- **Technical value:** PoDoFo delivers innovation to the PDF technical world
- **Cross community value:** Enhancements related to RELAX NG parsing integrated into libXML2



This presentation was sponsored by

DDDI EURONOVATE

The technology vendor for Trust Solutions



PDF Days Europe 2025

euronovategroup.com

Questions (to and from the audience)

- Why there was no public RNG grammar for the PDF XMP packets until now?
- Answering to anything (I try), also on PoDoFo...

Surf now to this ▼ presentation ▼



