# The secrets of a file

Jukka Zitting | contributor to Tika, PDFBox, POI, Commons, TagSoup, junrar, etc.

Thursday, November 8, 12

# Agenda

- Type detection

- Metadata extraction

- Future work: XMP

Thursday, November 8, 12

# Type detection

Thursday, November 8, 12

# Demo

Thursday, November 8, 12

# File extensions

Thursday, November 8, 12

# Source code: one standard extension

```xml
<mime-type type="text/x-java-source">
  <_comment>Java source code</_comment>
  <alias type="text/x-java" />
  <glob pattern="*.java"/>
  <sub-class-of type="text/plain"/>
</mime-type>
```

Thursday, November 8, 12

# Source code: multiple commonly used extensions

```xml
<mime-type type="text/x-c++src">
  <_comment>C++ source code</_comment>
  <glob pattern="*.cpp"/>
  <glob pattern="*.cxx"/>
  <glob pattern="*.cc"/>
  <glob pattern="*.C"/>
  <glob pattern="*.c++"/>
  <glob pattern="*.CPP"/>
  <sub-class-of type="text/plain"/>
</mime-type>
```

Thursday, November 8, 12

```
<mime-type type="text/x-prolog">
 <_comment>Prolog source</_comment>
 <glob pattern="*.pro"/>
 <!-- <glob pattern="*.pl"/>
     - conflicts with text/x-perl -->
 <sub-class-of type="text/plain"/>
</mime-type>
```

Thursday, November 8, 12

# Common first step also for other types of files

```
<mime-type type="image/x-raw-canon">
 <_comment>Canon raw image</_comment>
 <glob pattern="*.crw"/>
 <glob pattern="*.cr2"/>
</mime-type>
```

Thursday, November 8, 12

# No file name? Use magic!

Thursday, November 8, 12

# One standard byte pattern

```xml
<mime-type type="application/pdf">
  <alias type="application/x-pdf"/>
  <acronym>PDF</acronym>
  <_comment>Portable Doc...</_comment>
  <magic priority="50">
    <match value="%PDF-"
             type="string" offset="0"/>
  </magic>
  <glob pattern="*.pdf"/>
</mime-type>
```

Thursday, November 8, 12

# Alternative byte patterns

```xml
<mime-type type="image/gif">
 <acronym>GIF</acronym>
 <_comment>Graphics Inter...</_comment>
 <magic priority="50">
  <match value="GIF87a"
          type="string" offset="0"/>
  <match value="GIF89a"
          type="string" offset="0"/>
 </magic>
 <glob pattern="*.gif"/>
</mime-type>
```

Thursday, November 8, 12

# Odd cases: MS Word

```
<mime-type type="application/msword">
  <alias type="application/vnd.ms-word"/>
  <magic priority="50">
    <match value="Microsoft\ Word\ 6.0..."
              type="string" offset="2080"/>
    <match value="Documento\ Microsoft..."
              type="string" offset="2080"/>
    <match value="MSWordDoc"
              type="string" offset="2112"/>
    <match value="0x31be0000"
              type="big32" offset="0"/>
```

# Odd cases: HTML

```
<mime-type type="text/html">
  <magic priority="40">
    <match value="&lt;!DOCTYPE HTML"
            type="string" offset="0:64"/>
    <match value="&lt;HTML"
            type="string" offset="0:64"/>
    <match value="&lt;HEAD"
            type="string" offset="0:64"/>
    <match value="&lt;TITLE"
            type="string" offset="0:64"/>
    <match value="&lt;BODY"
```

Thursday, November 8, 12

# Container formats

Thursday, November 8, 12

# XML formats

```
<mime-type type="application/xhtml+xml">
  <magic priority="50">
    <match value="&lt;html xmlns="
            type="string" offset="0:8192"/>
  </magic>
  <root-XML namespaceURI=
            "http://www.w3.org/1999/xhtml"
            localName="html"/>
  <glob pattern="*.xhtml"/>
  <glob pattern="*.xht"/>
</mime-type>
```

Thursday, November 8, 12

# ZIP archives

```xml
<mime-type type="application/
        vnd.oasis.opendocument.spreadsheet">
 <magic>
  <match type="string" offset="0" value="PK">
   <match type="string" offset="30"
     value="mimetypeapplication/
       vnd.oasis.opendocument.spreadsheet"/>
  </match>
 </magic>
 <glob pattern="*.ods"/>
</mime-type>
```

Thursday, November 8, 12

# Custom cases

Thursday, November 8, 12

# Custom cases

- Zip archives
  - Parse the container and look at contained file names
- Old MS Office formats
  - Parse the container and look at contained resources
- Plain text
  - Really tricky, see below...

Thursday, November 8, 12

# Detecting plain text

- UTF BOMs
- Control characters
- Line endings
- Byte histogram
- Not foolproof, but quite reliable in practice

Thursday, November 8, 12

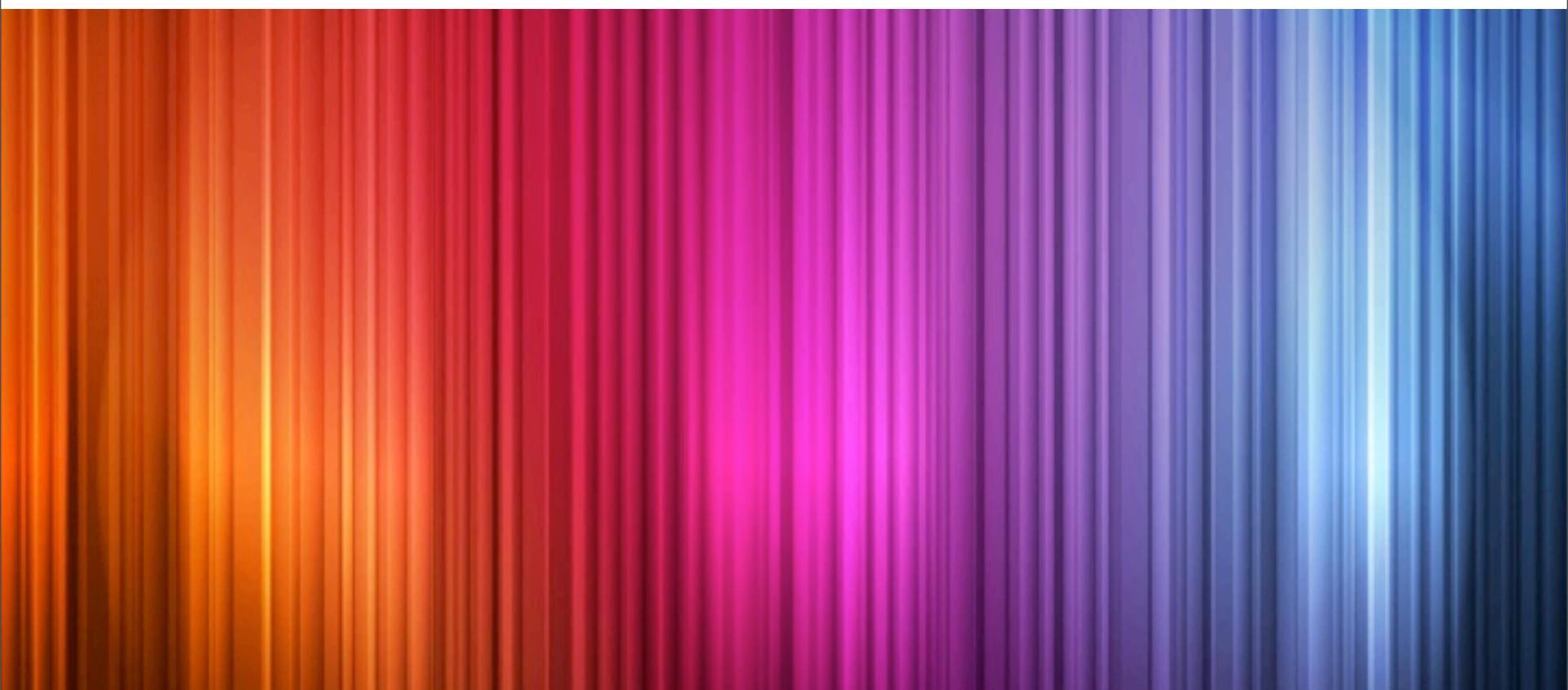# Composite approach

Thursday, November 8, 12

# The ultimate type detector

- Custom cases

- Magic byte patterns

- Container types

- File name hints

- Content type hints

- application/octet-stream


- Less reliable detectors can only add detail to more reliable results

Thursday, November 8, 12

# Metadata extraction

Thursday, November 8, 12

# Demo

Thursday, November 8, 12

# Types of metadata

Thursday, November 8, 12

# Dublin core

- Basic information
  - dc:title
  - dc:creator
  - dc:date
  - dc:format
- Driven originally by (scientific) libraries
- http://dublincore.org/

Thursday, November 8, 12

# Exif

- Image metadata
  - tiff:ImageLength
  - tiff:ImageWidth
  - tiff:BitsPerSample
  - …
- Useful also for non-TIFF/JPEG image formats

Thursday, November 8, 12

# Other schemas

- International Press Telecommunications Council (IPTC)
- Various XMP-related schemas

Thursday, November 8, 12

# Future work

Thursday, November 8, 12

# XMP

Thursday, November 8, 12