# Session:
# Characterisation of Digital Content

**Digital Preservation – The Planets Way**
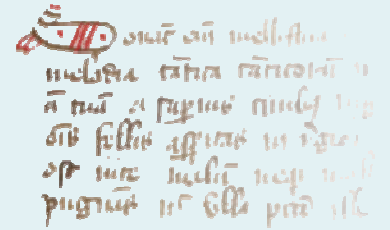**Sofia, 16 – 18 September 2009**

**Volker Heydegger and Jan Schnasse**

# Overview

❑ Part 1: Characterising Digital Content: The eXtensible Characterisation Languages

❑ Part 2: Demonstration of XCL Tools:
Evaluation of Format Conversion

# Characterising Digital Content: The eXtensible Characterisation Languages

**Digital Preservation – The Planets Way**
**Sofia, 16 – 18 September 2009**

**Volker Heydegger**

# Overview

- Characterisation: Why and What
- About File Formats
- XCL: Goals
- XCL: Architecture
- XCL by Example

# Characterisation

## Why characterisation?

"Characterisation is an essential precursor to preservation. It provides the information required to make preservation planning decisions about digital objects, and to validate the results of preservation actions. "

(A. Brown: Developing Practical Approaches to Active Preservation, IJDC, 2007)

# Characterisation

## Why characterisation?

"Characterisation is an essential precursor to preservation. It provides the information required to make preservation planning decisions about digital objects, and to validate the results of preservation actions. "

(A. Brown: Developing Practical Approaches to Active Preservation, IJDC, 2007)
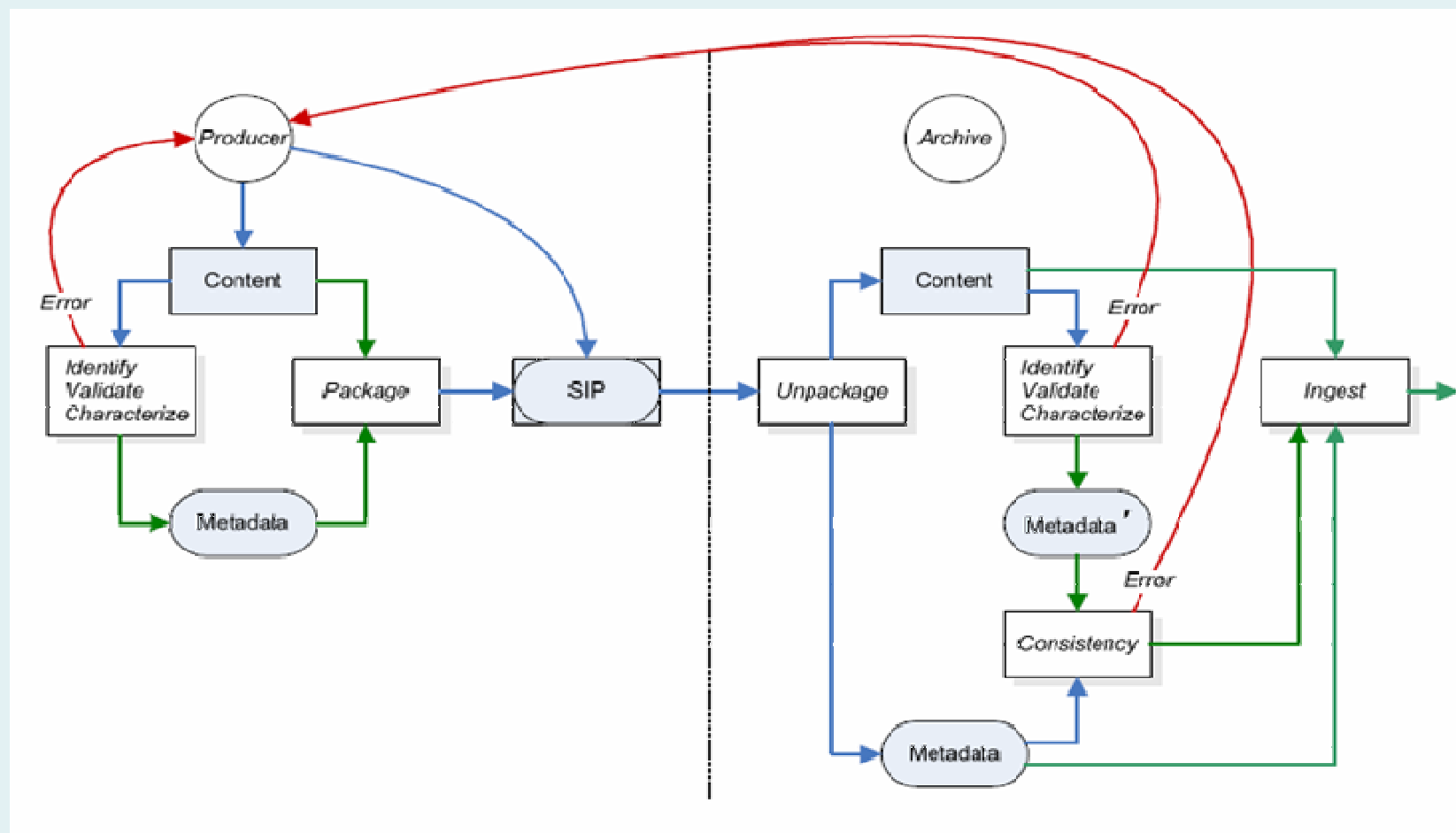
# Characterisation

## Why characterisation?

"Characterisation is an essential precursor to preservation. It provides the information required to make preservation planning decisions about digital objects, and to validate the results of preservation actions. "

(A. Brown: Developing Practical Approaches to Active Preservation, IJDC, 2007)

# Why characterisation?

# Characterisation

## What is subject to characterisation?

"One essential process in digital preservation is to perform format characterization to extract technical metadata associated with each digital object in the preservation archival collection. The technical metadata are important attributes for understanding and managing the digital archival collections, especially for format monitoring and researching format transformation procedures."

(C.C.H. Chou: Format Identification, Validation, Characterization and Transformation in DAITSS, [?2007])

# Characterisation

What is subject to characterisation?

"One essential process in digital preservation is to perform format characterization to extract technical metadata associated with each digital object in the preservation archival collection. The technical metadata are important attributes for understanding and managing the digital archival collections, especially for format monitoring and researching format transformation procedures."

(C.C.H. Chou: Format Identification, Validation, Characterization and Transformation in DAITSS, [?2007])

# Characterisation

## What is subject to characterisation?

"One essential process in digital preservation is to perform format characterization to extract technical metadata associated with each digital object in the preservation archival collection. The technical metadata are important attributes for understanding and managing the digital archival collections, especially for format monitoring and researching format transformation procedures."
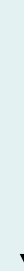
(C.C.H. Chou: Format Identification, Validation, Characterization and Transformation in DAITSS, [?2007])

# About File Formats

## What is a format?

0111001100011101000011010...

❑ On a very basic level (storage level) digital content is nothing but binary data

❑ On the software level, digital content is stored as formatted data, i.e. as *meaningful* sequences of bytes

→ (*File) Format*

❑ On the most human-perceivable level it appears in a rendered form

# How many file formats?

- PRONOM: ~ 550

- www.wotsit.org: ~ 900

- www.fileformat.info:  567

- www.fileinfo.com: > 3000   (file extensions)

# How many file formats can we find in institutions?

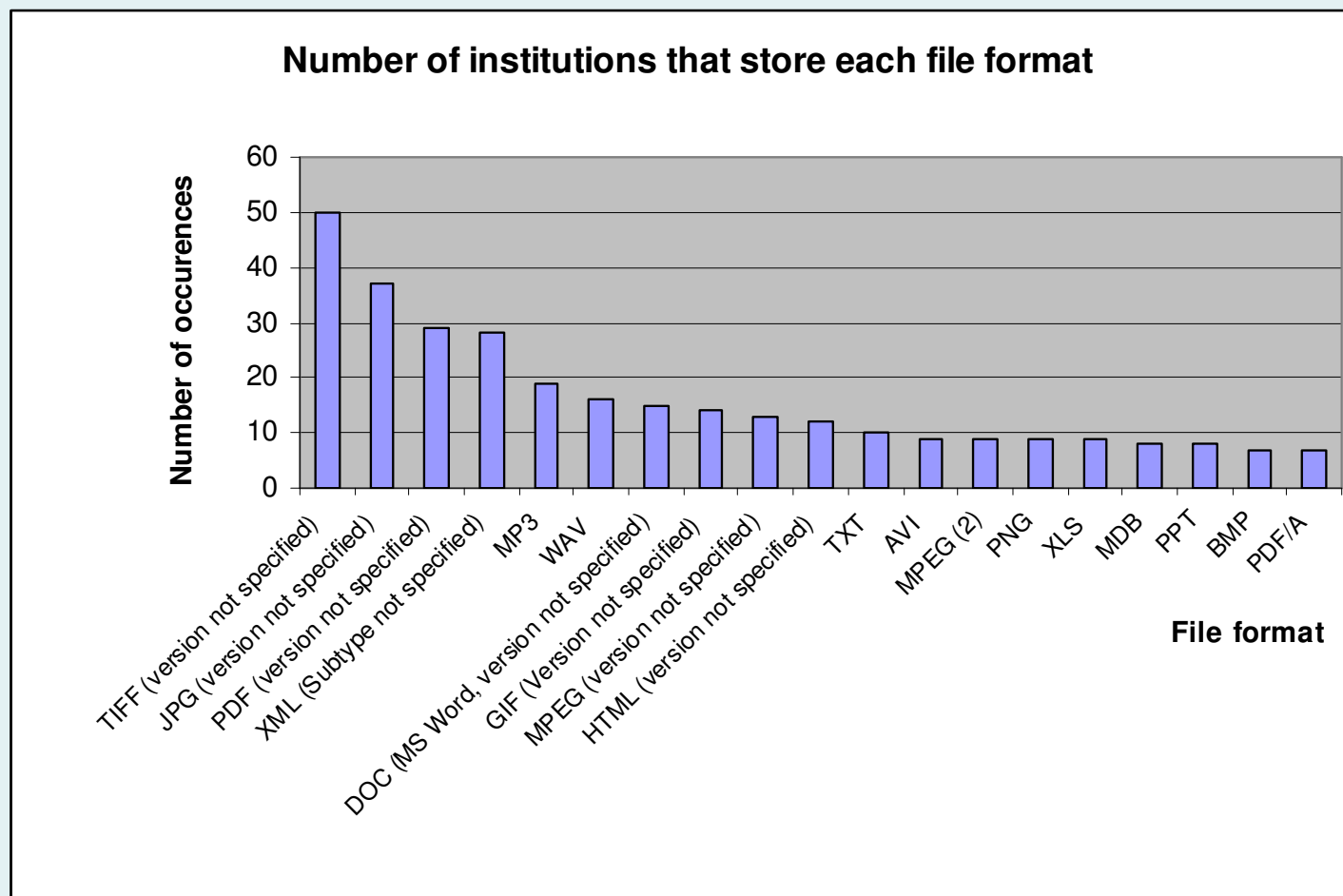Planets internal study: "Gap analysis in tool provision"

- 76 institutions from 13 countries

- 137 different file formats (124 excl. versions)

# How many file formats are used more often?

**Number of institutions that store each file format**



Source: Planets internal report:
Gap analysis in tool provision (third version).

# Suitability of formats for preservation (1)

| High confidence | Medium confidence | Low confidence |
|---|---|---|
| ❖TIFF (uncompressed)<br>❖ PNG (*.png) | ❖ BMP (*.bmp)<br>❖ JPEG/JFIF (*.jpg)<br>❖JPEG2000 (prefer lossless or uncompressed) (*.jp2)<br>❖TIFF (compressed)<br>❖GIF (*.gif) | ❖MrSID (*.sid)<br>❖TIFF (in Planar format)<br>❖FlashPix (*.fpx)<br>❖PhotoShop (*.psd)<br>❖All other raster image formats not listed here |

Source: http://www.fcla.edu/digitalArchive/ pdfs/recFormats.pdf

CEI CENTRAL EUROPEAN INITIATIVE

planets

# Suitability of formats for preservation (2)

| High confidence | Medium confidence | Low confidence |
|---|---|---|
| ❖ Plain text (encoding: ISO8859-1 - 9, UTF-8, UTF-16 with BOM)<br>❖ XML (includes XSD/XSL/XHTML, etc.; with included or accessible schema and character encoding explicitly specified)<br>❖ PDF/A-1 (ISO 19005-1) | ❖ Cascading Style Sheets (*.css)<br>❖ DTD (*.dtd)<br>❖ PDF (*.pdf) (embedded fonts)<br>❖ Rich Text Format 1.x (*.rtf)<br>❖ HTML 4.x (include a DOCTYPE declaration)<br>❖ SGML (*.sgml)<br>❖ Open Office (*.sxw/*.odt)<br>❖ Office Open XML (*.docx) | ❖PDF (*.pdf) (encrypted)<br>❖ Microsoft Word (*.doc)<br>❖ WordPerfect (*.wpd)<br>❖ DVI (*.dvi)<br>❖ All other text formats not listed here |

CEI
CENTRAL EUROPEAN INITIATIVE

planets

# Suitability of formats for preservation (3)

| High confidence | Medium confidence | Low confidence |
|---|---|---|
| ❖AIFF (PCM) (*.aif, *.aiff)<br>❖ WAV (PCM) (*.wav) | ❖SUN Audio (uncompressed) (*.au)<br>❖Standard MIDI (*.mid, *.midi)<br>❖Ogg Vorbis (*.ogg)<br>❖Free Lossless Audio Codec (*.flac)<br>❖ Advance Audio Coding (*.mp4, *.m4a, *.aac)<br>❖ MP3 (MPEG-1/2, Layer 3)(*.mp3) | ❖AIFC (compressed) (*.aifc)<br>❖ NeXT SND (*.snd)<br>❖ RealNetworks 'Real Audio, (*.ra, *.rm, *.ram)<br>❖ Windows Media Audio<br>❖(*.wma)<br>❖WAV (compressed) (*.wav)<br>❖All other audio formats not listed here |

# Criteria for suitability

- Openess
- Adoption
- Complexity
- Technical protection mechanism
- Self-documentation
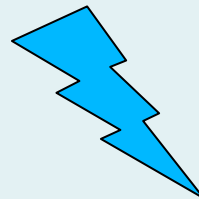- Robustness
- Dependencies

(J. Rog, C. van Wijk: Evaluating File Formats for Long-term Preservation, iPres 2007)
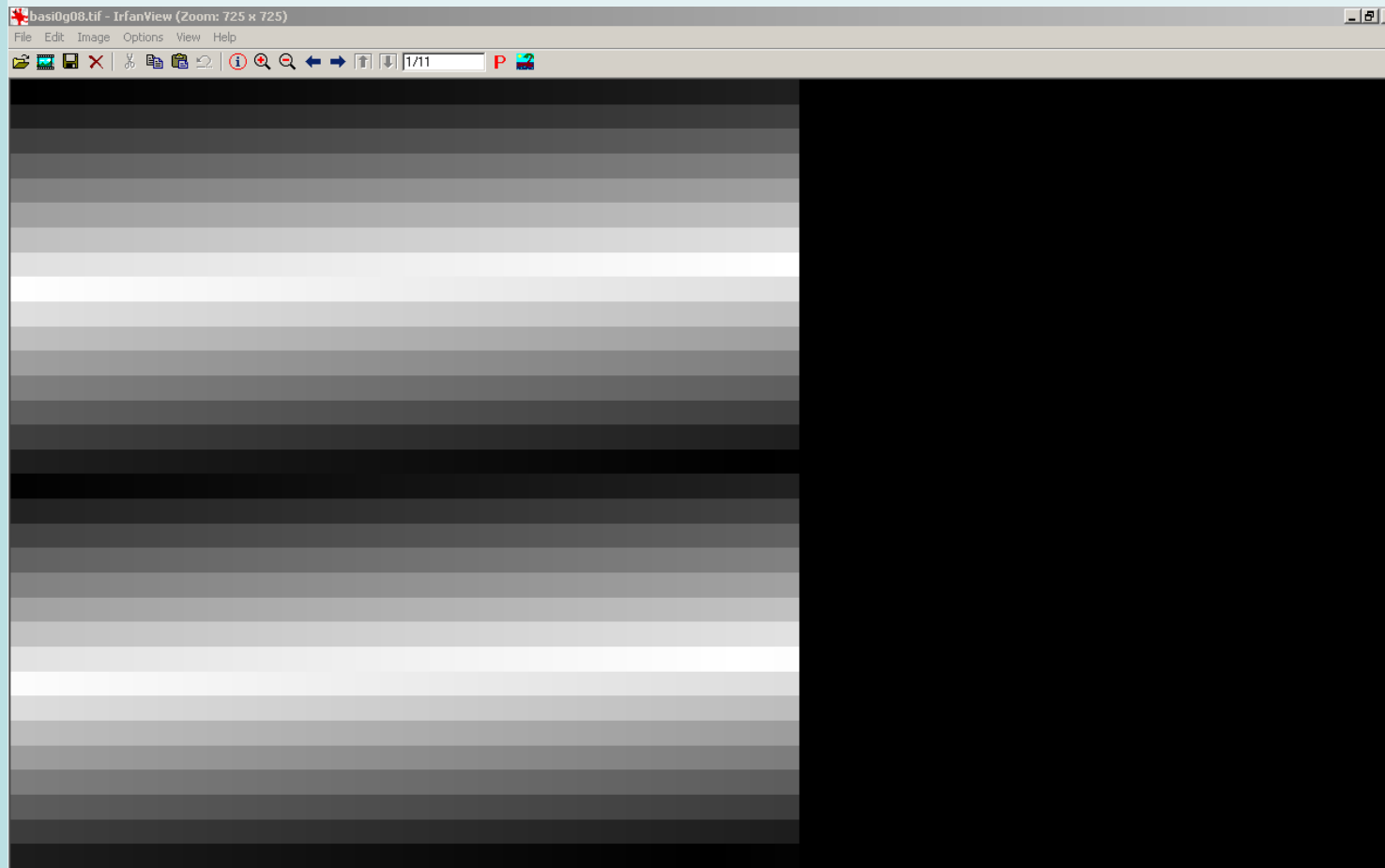
# Robustness of Formats

*Robustness*

    ::= resilience of file formats against bit-stream corruption
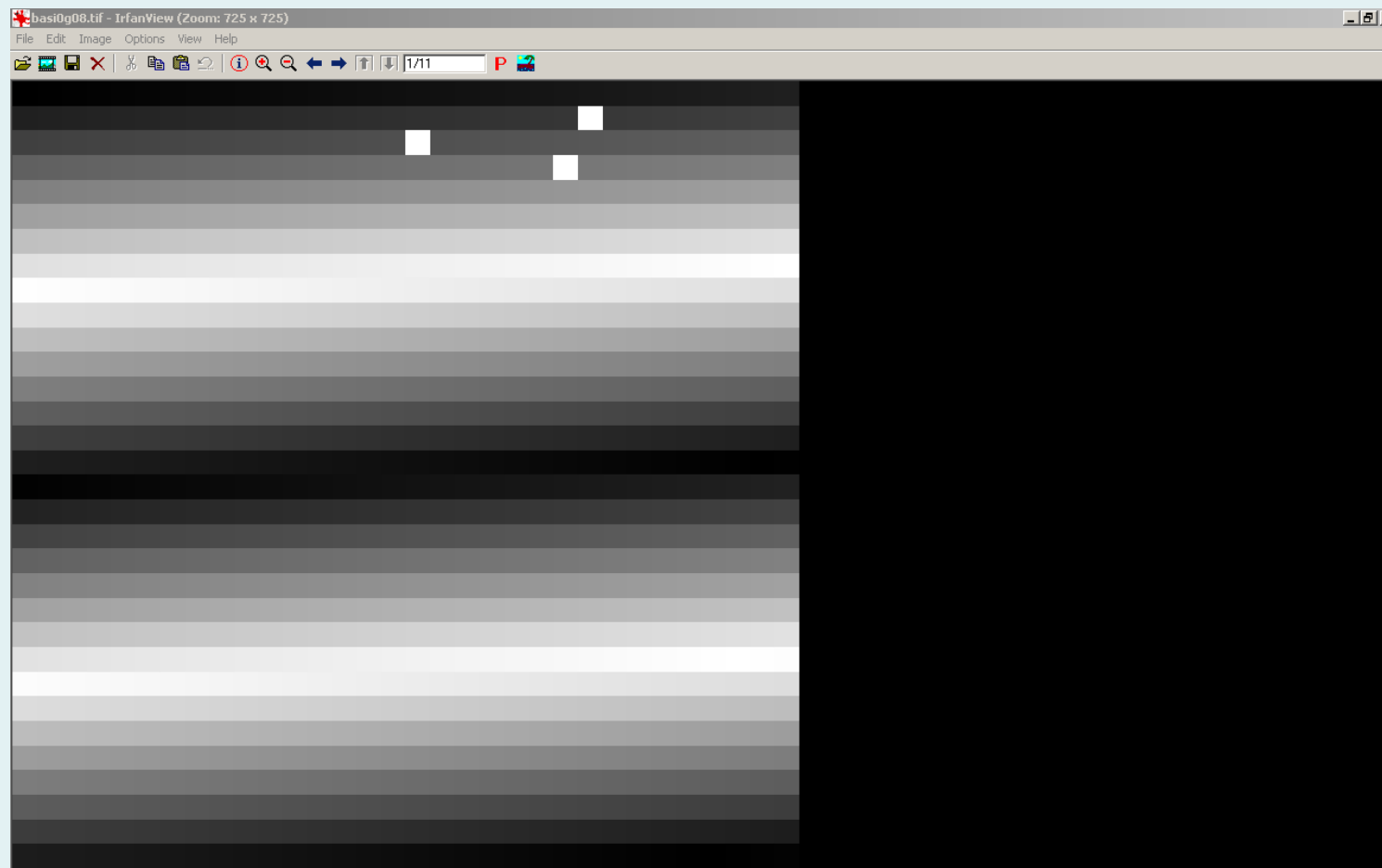
# What happens if data is corrupted in files?



**Testimage: Tiff, greyscale, 32x32 pixel, 8 bit per pixel**

**First 224 bytes of testfile**

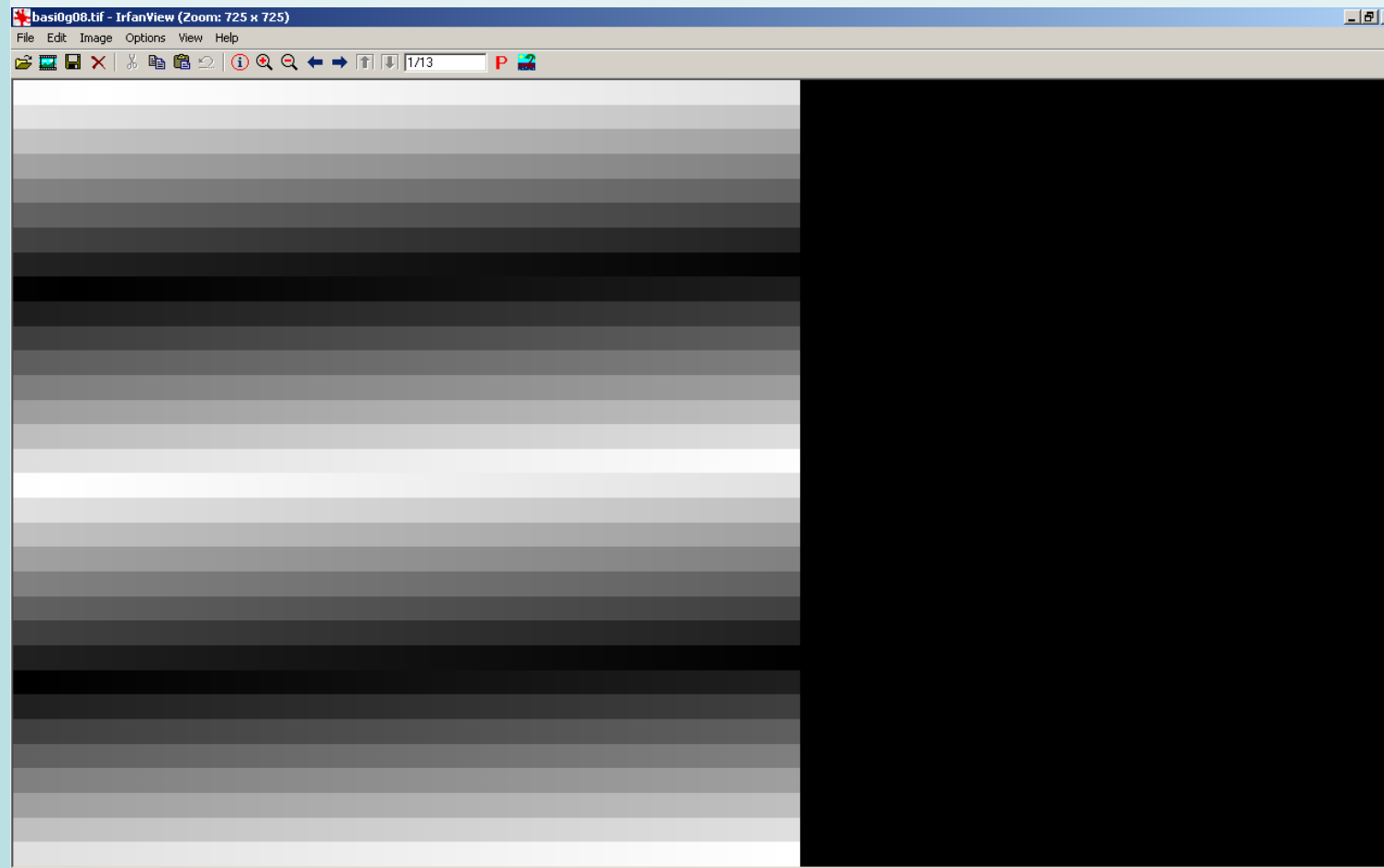# Information loss: 1 byte data = = 1 Pixel

```
0x3F0:  14 13 12 11 10 0F 0E 0D 0C 0B 0A 09 08 07 06 05
0x400:  04 03 02 01 00 01 02 03 10 00 00 01 03 00 01 00
0x410:  00 00 20 00 00 00 01 01 03 00 01 00 00 00 20 00
0x420:  00 00 02 01 03 00 01 00 00 00 08 00 00 00 03 01
0x430:  03 00 01 00 00 00 01 00 00 00 06 01 03 00 01 00
0x440:  00 00 01 00 00 00 11 01 04 00 01 00 00 00 08 00
0x450:  00 00 12 01 03 00 01 00 00 00 03 00 00 00 15 01
0x460:  03 00 01 00 00 00 01 00 00 00 16 01 03 00 01 00
0x470:  00 00 00 01 00 00 17 01 04 00 01 00 00 00 00 04
0x480:  00 00 1A 01 05 00 01 00 00 00 CE 04 00 00 1B 01
0x490:  05 00 01 00 00 00 D6 04 00 00 1C 01 03 00 01 00
0x4A0:  00 00 01 00 00 00 28 01 03 00 01 00 00 00 02 00
0x4B0:  00 00 31 01 02 00 0A 00 00 00 DE 04 00 00 40 01
0x4C0:  03 00 00 03 00 00 E8 04 00 00 00 00 00 00 00 00
```

00

**Part of the TIFF Image File Directory, Tag: Photometric Interpretation**

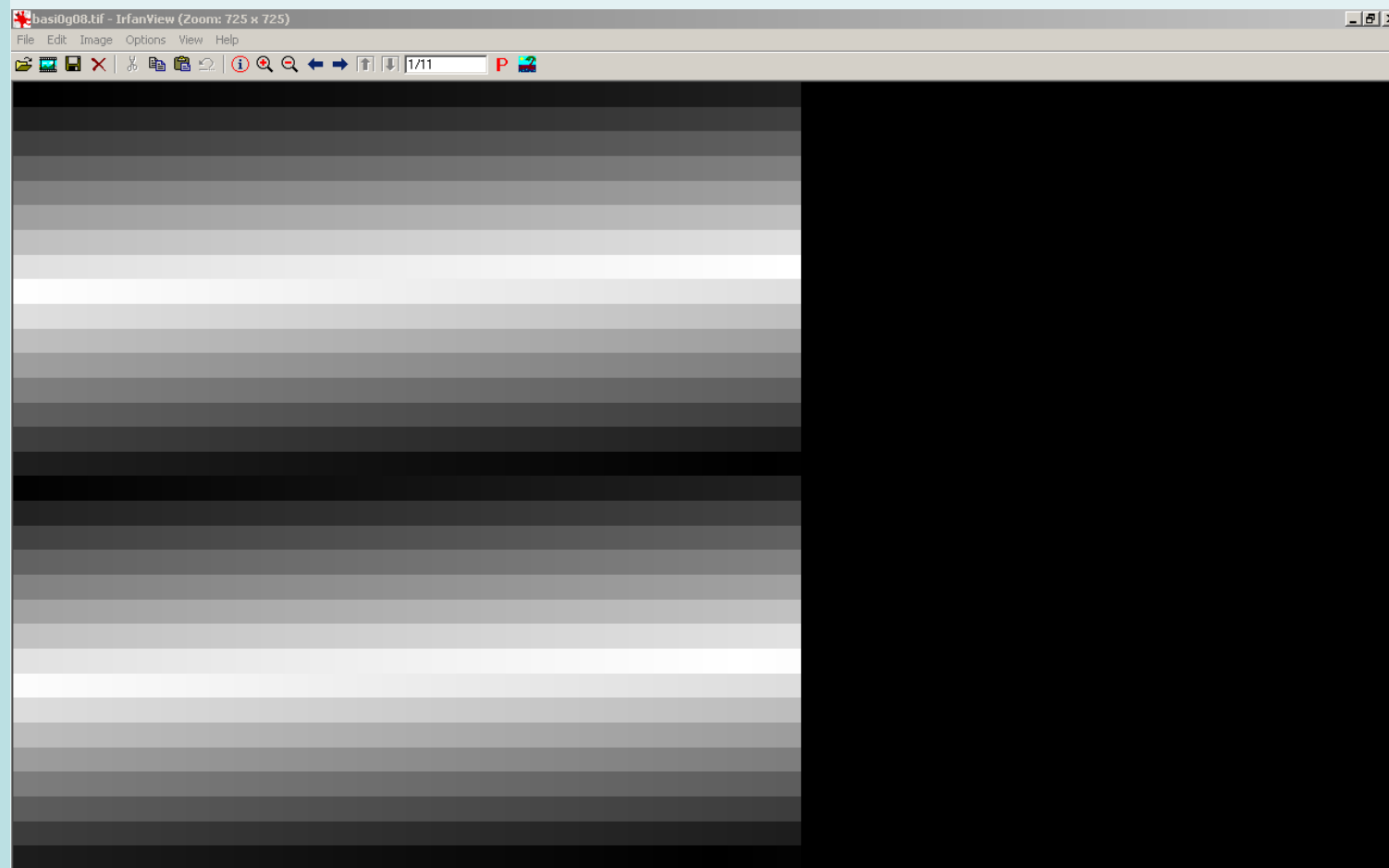# 1 bit changes == 100% information changed

Table 1: Results for $R_{Bt}$ (in percentage) for various file formats

| | 1 Byte | 0.01 | 0.1% | 1.0% |
|---|---|---|---|---|
| **TIFF** | | | | |
| uncompressed | 0.00 (0.00063) | 0.56 | 6.64 | 48.83 |
| JPEG compressed, ratio 1:2.60 (62%) | 2.14 (0.00166) | 13.03 | - | - |
| JPEG compressed, ratio 1:10.72 (90%) | 2.44 (0.00505) | 13.32 | - | - |
| LZW compressed, ratio 1:1.01 (2%) | 1.37 (0.00064) | 18.79 | 77.95 | 99.34 |
| ZIP compressed, ratio 1:1.28 (22%) | 27.12 (0.00081) | 84.92 | 98.47 | - |
| **PNG** | | | | |
| ZLIB compressed, unfiltered | 18.21 (0.00074) | 79.15 | 97.63 | - |
| ZLIB compressed, filtered | 25.05 (0.00085) | 81.83 | 98.08 | - |
| **BMP (windows)** | | | | |
| uncompressed | 0.00 (0.00063) | 0.14 | 1.92 | 15.29 |
| **JP2** | | | | |
| lossless, ratio 1:1.36 (27%) | 17.53 (0.00086) | 76.22 | 94.29 | - |
| lossy, ratio 1:7.42 (87%) | 33.31 (0.00166) | 51.86 | 95.03 | - |
| lossy, ratio 1:2.64 (62%) | 22.61 (0.00468) | 72.93 | 95.62 | - |

V.Heydegger: Analysing the Impact of File Formats on Data Integrity, Archiving 2008

# Categories of characteristics

What is subject to characterisation?

"One essential process in digital preservation is to perform format characterization to extract technical metadata associated with each digital object in the preservation archival collection. The technical metadata are important attributes for understanding and managing the digital archival collections, especially for format monitoring and researching format transformation procedures."

(C.C.H. Chou: Format Identification, Validation, Characterization and Transformation in DAITSS, [?2007])

# Non-technical characteristics ("associated metadata")

What's the name of the object?

Which software created the object?

Who holds the intellectual rights for the object?

When was the object modified for the last time?

Which collection does the object belong to?

Where is the object located in our repository?
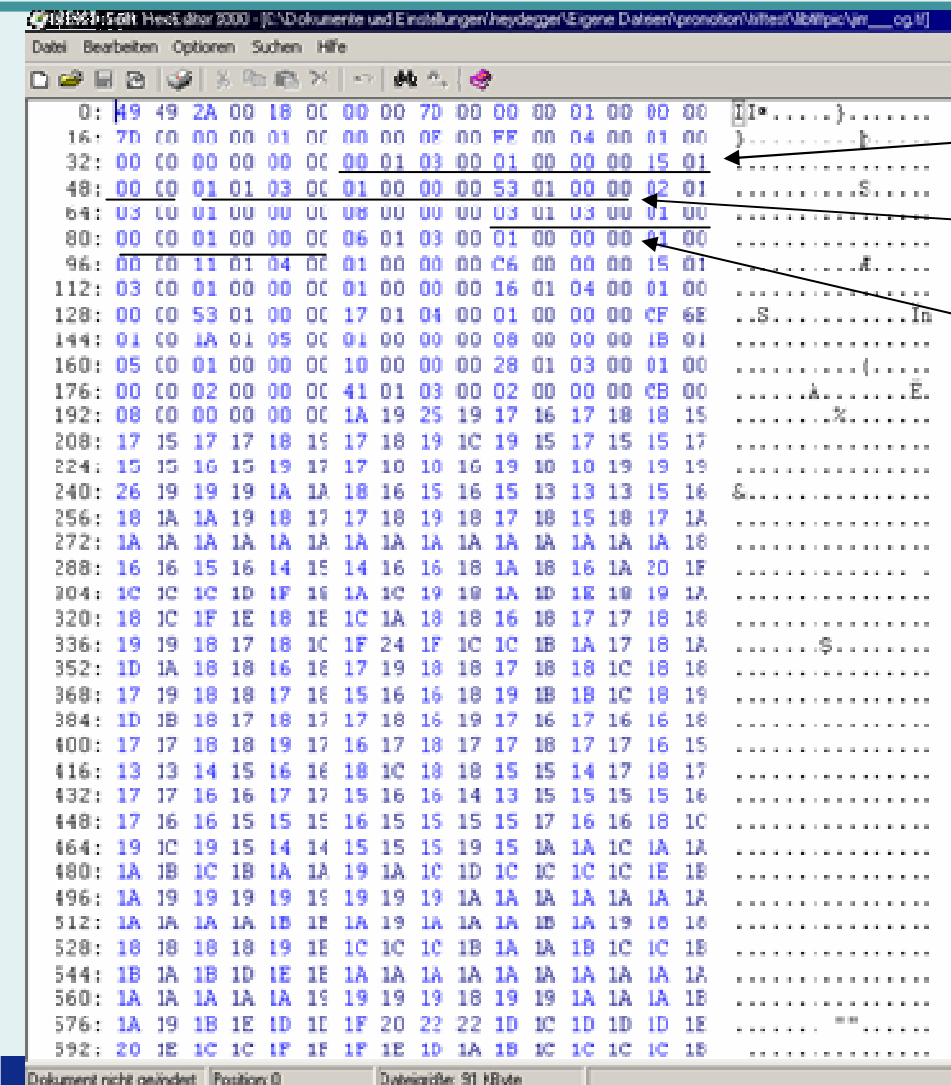
…

# Technical characteristics



Image width: 277

Image length: 339

Compression: uncompressed

## ImageLength

The number of rows of pixels in the image.

Tag = 257 (101.H)

Type = SHORT or LONG

N = 1

No default. See also ImageWidth.

## ImageWidth

The number of columns in the image, i.e., the number of pixels per row.

Tag = 256 (100.H)

Type = SHORT or LONG

N = 1

No default. See also ImageLength.

# Categories of characteristics

- Significant characteristics:

= Those properties which are essential for keeping the integrity of the object

→ Significant properties are always of technical nature

# Lessons learnt so far

- Characterisation is an essential part within an overall preservation framework.

- File Format is the central concept for representation of digital content.

- A Format describes the characteristics of objects.

- There is a huge amount of formats but only a couple of them are actually suitable for preservation.

# XCL: Goals

- Support preservation planning framework

- Support a specific preservation action task: Evaluation of file format conversion

- Develop a more abstract model for extraction of characteristics (syn. properties) from files

- Develop tools which use this model in order to enable characterisation in an efficiently, i.e. in an automated way

# XCL: Goals

- In practice:

   -   Develop an „eXtensible Characterisation Definition Language" (XCDL), able to describe the content of digital objects (=1 + n more files), processible by a software tool for further analysis.

   -   Develop an „eXtensible Characterisation Extraction Language" (XCEL), able to describe any machine readable format in a formal language,  processible by a software tool for extraction of content as XCDL.

# XCL: Goals

- Support preservation planning framework

- Support a specific preservation action task: Evaluation of file format conversion

- Develop a more abstract model for extraction of characteristics (syn. properties) from files

- Develop tools which use this model in order to enable characterisation in an efficiently, i.e. in an automated way

# Why automate?

Assumption:

Preservation is only feasible, if the content of two digital objects can be compared without human intervention.

# Why automate?

1 million objects: use *five minutes* for each

== 416 666.7    hours

== 52 803.4     8-hour days for a Human

# Why automate?

1 million objects: use *one second* for each.

== 16666.7 minutes == 277.8 hours

== 11.57 working days of a computer

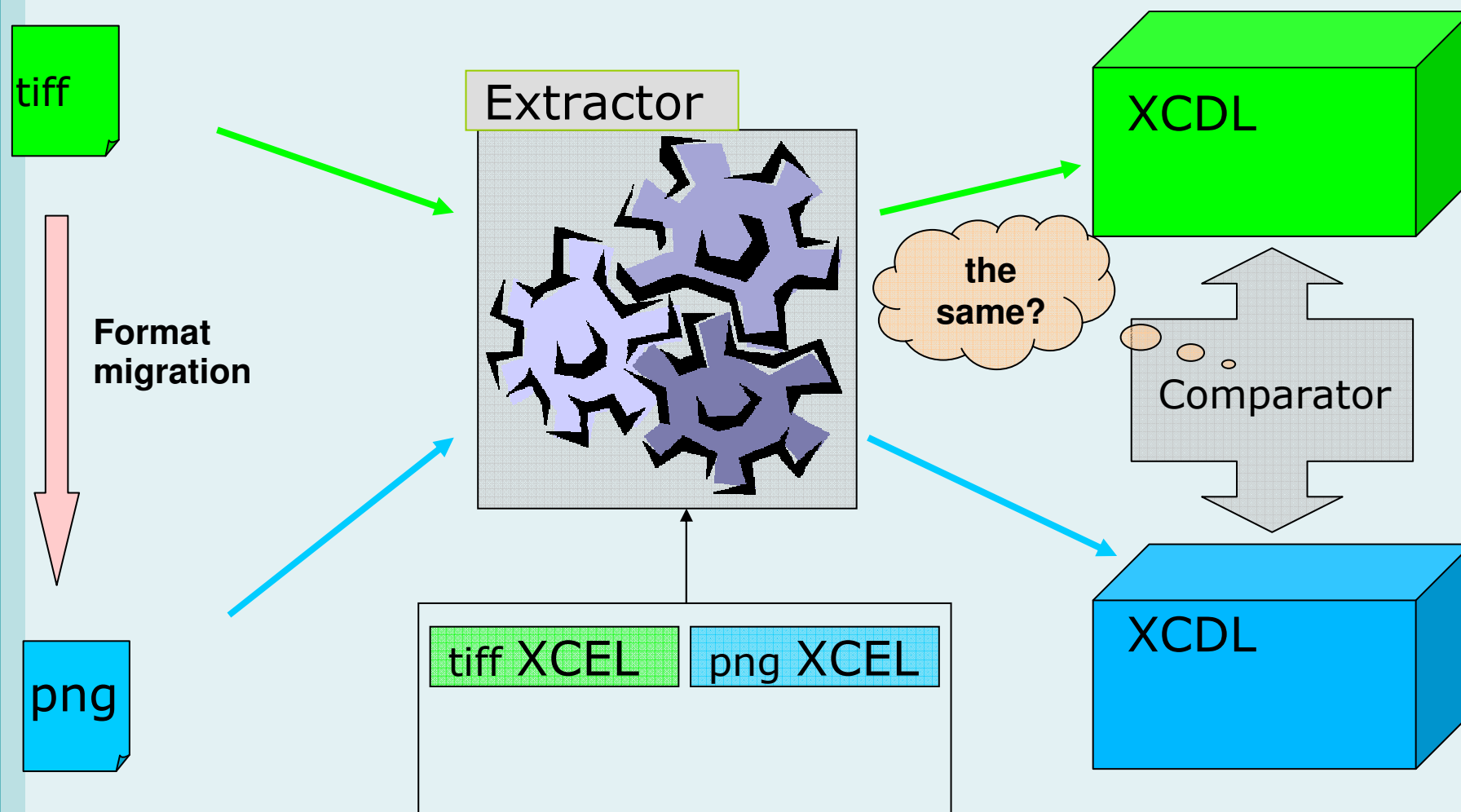== 34.7 8-hour days for a Human

== 7 working weeks

# Why automate?

# XCL: Goals

- Support preservation planning framework

- Support a specific preservation action task: Evaluation of file format conversion

- Develop a more abstract model for extraction of characteristics (syn. properties) from files

- Develop tools which use this model in order to enable characterisation in an efficiently, i.e. in an automated way

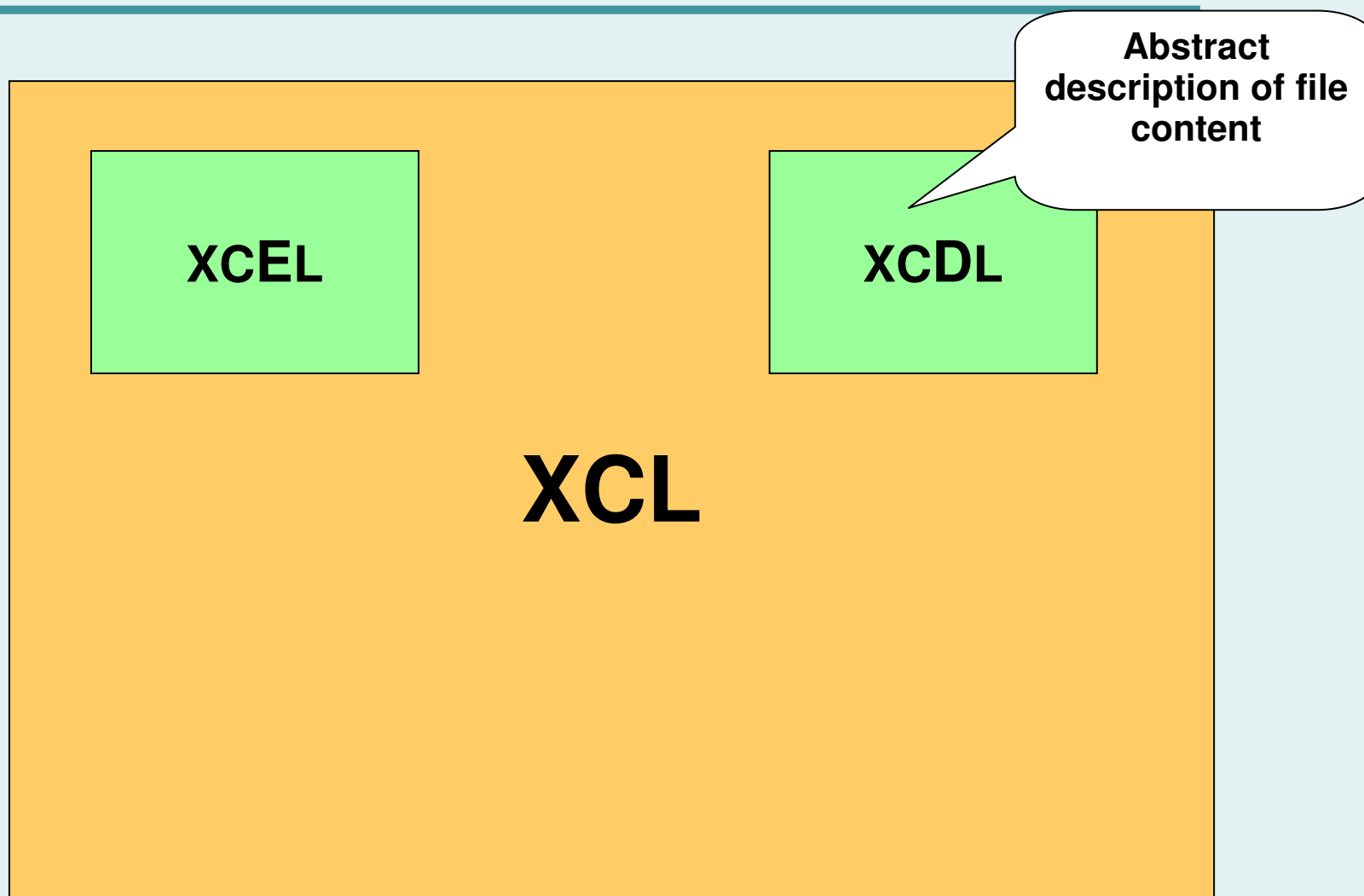# XCL: Main application: Evaluation of format conversion

# XCL: Architecture

**XCL**

# The Ontology

# XML as backbone language

# eXtensible Characterisation Extraction Language (XCEL)

➢ Describing how properties of digital objects are stored

➢ File format specification tagged in XML, according to the XCEL language definitions

➢ Interpretable through an XCEL interpreter (Extractor), able to extract characteristics

# XCEL: Global Architecture

**XCEL Description**

**Preprocessing** ← Configuration tasks, affecting the behaviour of the XCEL interpreter

**Format description** ← Description of the structure of the object

**Templates** ← Description of recuring structures

**Postprocessing** ← Actions on the result of the format description processing

# XCEL: Basic Structuring Elements

**There are just a few elements sufficient enough to describe a file format:**

processing

nonValidValues

valueInterpretation

param

item

valueLabel

value

symbol

# eXtensible Characterisation Definition Language (XCDL)

- Describes the content of a file /set of files in an abstract way.

- Designed for decription of the content of *any* file format.

- Designed as a means to describe only parts *or* all of the content.

# XCDL: Basic Structuring Elements

**Again, there are just a few elements sufficient enough to describe the content of a digital object:**

object

type

dataRef

normData

property

valueSet

value

labValue

propertySet

# Benefits of the XCL approach

- XCL is a generic solution, uses an abstract model, provides a unique vocabulary

  →Extensible: XCL is based on XML

  →XCEL provides a means for description of any file format

  → XCDL is a language with which all sort of content can be  expressed

# XCL by Example



Image width: 277

Image length: 339

## ImageLength

The number of rows of pixels in the image.

Tag = 257 (101.H)

Type = SHORT or LONG

N = 1

No default. See also ImageWidth.

## ImageWidth

The number of columns in the image, i.e., the number of pixels per row.

Tag = 256 (100.H)

Type = SHORT or LONG

N = 1

No default. See also ImageLength.

# XCEL representation

```xml
<!-- Tag 256: ImageWidth (XCL: imageWidth) -->
  <item xsi:type="structuringItem" identifier="IFDE_256"
   optional="true">
   <symbol interpretation="uint16" length="2" value="256"/>
   <item xsi:type="structuringItem" order="choice">
    <item xsi:type="structuringItem" order="sequence">
     <!– Data type (value ,3' means uint16)-->
     <symbol interpretation="uint16" length="2"  value="3"/>
     <!– number of values (N)->
     <symbol interpretation="uint32" length="4" value="1"/>
     <!-- the value and name of property -->
     <symbol interpretation="uint16" length="2"
     name="imageWidth"/>
     <!-- wasted space-->
     <symbol interpretation="uint16" length="2"/>
     […]
    </item>
   </item>
  </item>
```
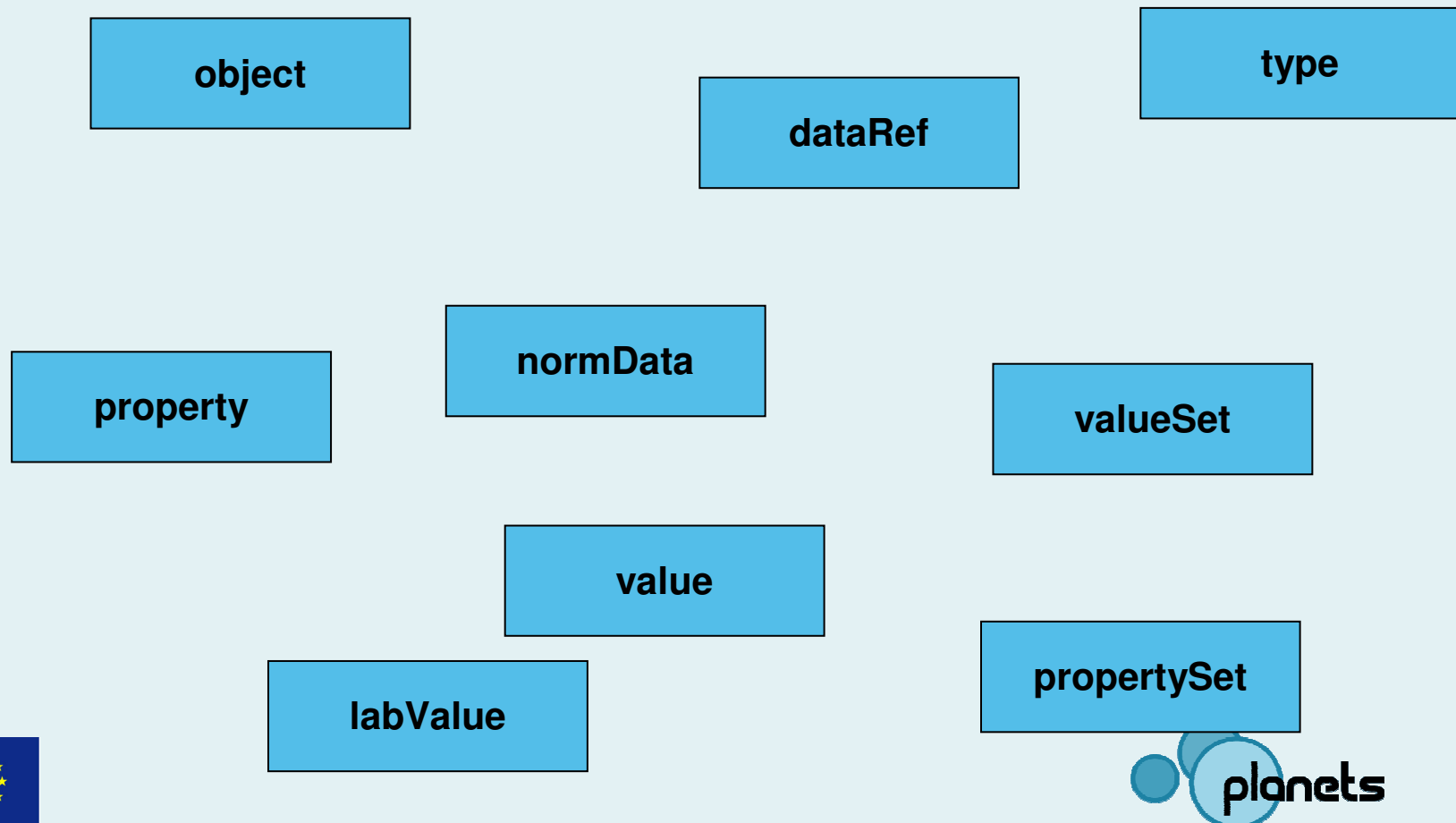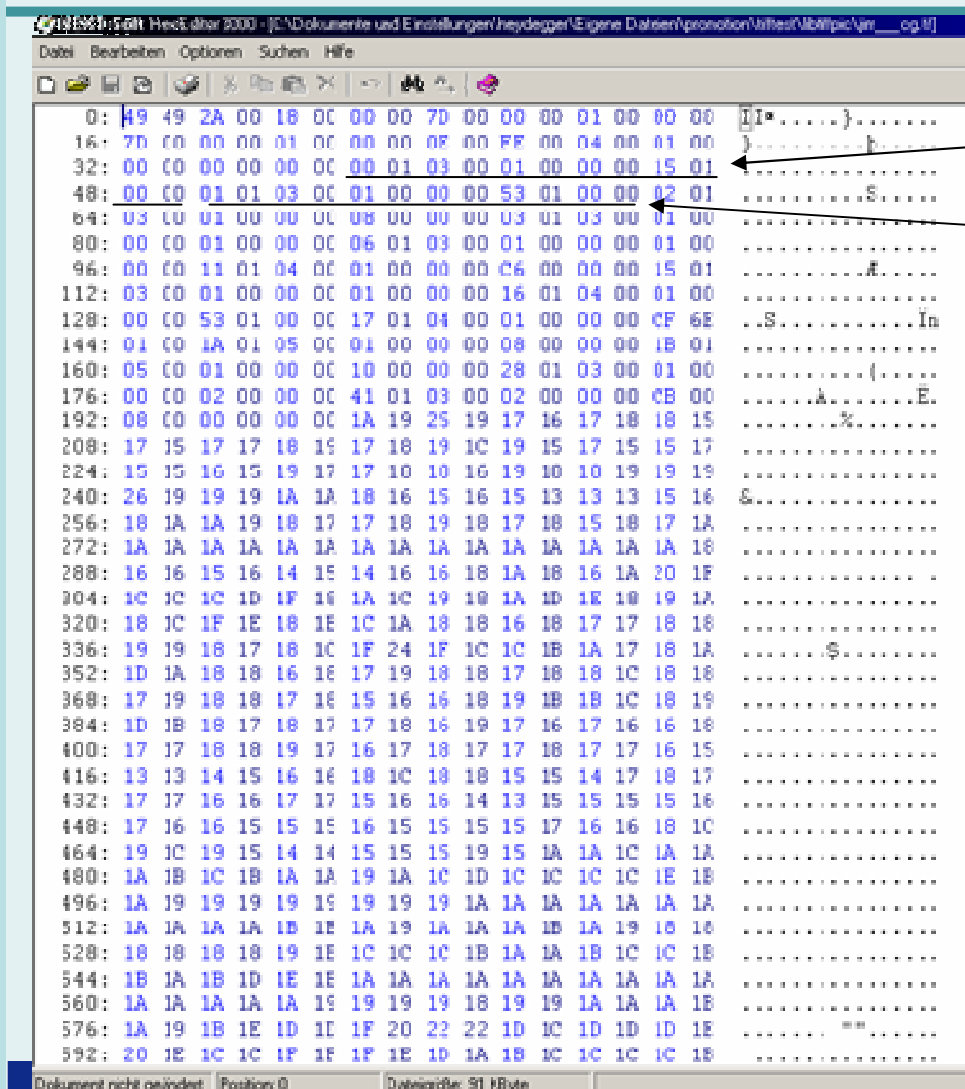
*ImageWidth*

The number of columns in the image, i.e., the number of pixels per row.

Tag     = 256 (100.H)

Type   = SHORT or LONG

N       = 1

No default. See also ImageLength.

# XCEL representation

```
<!-- Tag 256: ImageWidth (XCL: imageWidth) -->
 <item xsi:type="structuringItem" identifier="IFDE_256"
  optional="true">
 <symbol interpretation="uint16" length="2" value="256"/>
 <item xsi:type="structuringItem" order="choice">
  <item xsi:type="structuringItem" order="sequence">
   <!– Data type (value ‚3' means uint16)-->
   <symbol interpretation="uint16" length="2"  value="3"/>
   <!– number of values (N)->
   <symbol interpretation="uint32" length="4" value="1"/>
   <!-- the value and name of property -->
   <symbol interpretation="uint16" length="2"
   name="imageWidth"/>
   <!-- wasted space-->
   <symbol interpretation="uint16" length="2"/>
   […]
  </item>
 </item>
</item>
```

*ImageWidth*

The number of columns in the image, i.e., the number of pixels per row.

Tag     = 256 (100.H)

Type   = SHORT or LONG

N       = 1

No default. See also ImageLength.

# XCEL representation

```
<!-- Tag 256: ImageWidth (XCL: imageWidth) -->
 <item xsi:type="structuringItem" identifier="IFDE_256"
  optional="true">
 <symbol interpretation="uint16" length="2" value="256"/>
 <item xsi:type="structuringItem" order="choice">
  <item xsi:type="structuringItem" order="sequence">
   <!– Data type (value ,3' means uint16)-->
   <symbol interpretation="uint16" length="2"  value="3"/>
   <!– number of values (N)->
   <symbol interpretation="uint32" length="4" value="1"/>
   <!-- the value and name of property -->
   <symbol interpretation="uint16" length="2"
   name="imageWidth"/>
   <!-- wasted space-->
   <symbol interpretation="uint16" length="2"/>
   […]
  </item>
 </item>
 </item>
```

*ImageWidth*

The number of columns in the image, i.e., the number of pixels per row.

Tag     = 256 (100.H)

Type    = SHORT or LONG

N       = 1

No default. See also ImageLength.

# XCDL representation

```
...
<property id="p5">
    <name id="id30" >imageWidth</name>
    <valueSet id="i_i1_s4" >
        <labValue>
            <val>277</val>
            <type>int</type>
        </labValue>
    </valueSet>
  </property>
...
```

**XCEL entry:**
```
<!-- the value and name of property -->
    <symbol interpretation="uint16" length="2"
    name="imageWidth"/>
```

# XCDL representation

```
...
<property id="p5">
    <name id="id30" >imageWidth</name>
    <valueSet id="i_i1_s4" >
        <labValue>
            <val>277</val>
            <type>int</type>
        </labValue>
    </valueSet>
  </property>
...
```

**XCEL entry:**
<!– Data type (value ,3' means uint16)-->
   <symbol interpretation="uint16"
    length="2"  value="3"/>

# XCDL representations can now be compared...

```
Measure name: equal

Id: 1
Explanation: Metric 'equal' is a simple comparison of two values (A, B) of any XCL data type on
equality.
Data type of input value: Any XCL data type
Data type of output value: XCL: boolean (true, false)

Example:
```

| Value for property X of XCDL1 (src) | Value for property X of XCDL2 (tar) |
|---|---|
| `<labValue>`<br>`    <val>32</val>`<br>`    <type>int</type>`<br>`</labValue>` | `<labValue>`<br>`    <val>32</val>`<br>`    <type>int</type>`<br>`</labValue>` |

```
copra output:

…
<property id="2" name="imageHeight" unit="pixel" state="complete">
    <metrics>
        <metric id="1" name="equal">
            <result state="ok">true</result>
        </metric>
    </metrics>
</property>
…
```

Thank you for your attention!

Any questions?