


Historisch
Kulturwissenschaftliche
Informationsverarbeitung

Session: **Characterisation of Digital Content**

Digital Preservation – The Planets Way
Bern, 17 – 19 November 2009

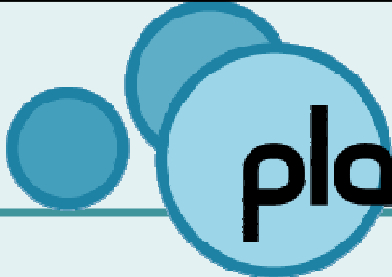
Volker Heydegger and Jan Schnasse



Overview

- ❑ Part 1: Characterising Digital Content: The eXtensible Characterisation Languages
- ❑ Part 2: Demonstration of XCL Tools: Evaluation of Format Conversion





planets


*Donat an indelible
nubila rāra rāra rāra
ā nū a pūpū nūpū
sū fūfū āfūfū nū fūfū
op nū nūfū nūfū nūfū
pūpū nū fūfū nūfū*

Historisch
Kulturwissenschaftliche
Informationsverarbeitung

Characterising Digital Content: The eXtensible Characterisation Languages

Digital Preservation – The Planets Way
Bern, 17 – 19 November 2009

Volker Heydegger



Overview

- ❑ Characterisation: General issues
- ❑ About File Formats
- ❑ XCL overview
- ❑ XCL by Example



Characterisation: General issues

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: General issues

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: General issues

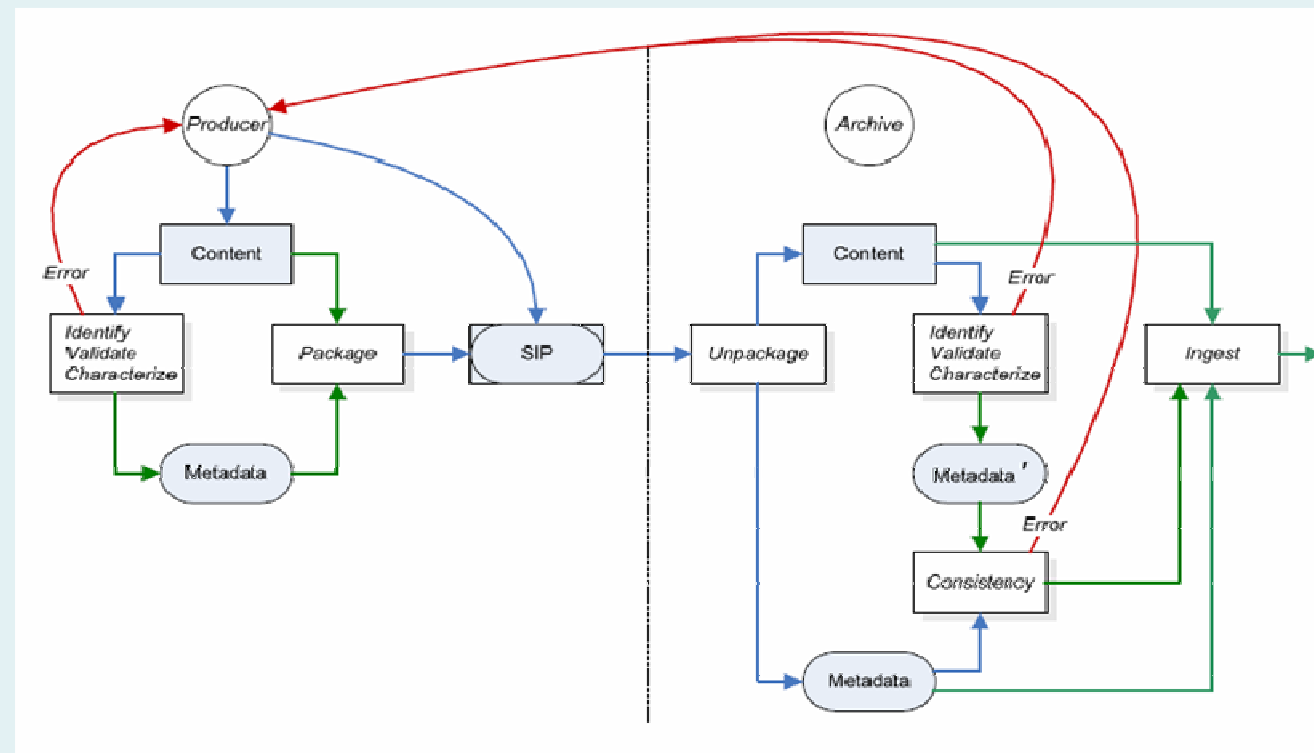
Why characterisation?

“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?



Source: S. Abrams: Automated Characterization in Preservation Workflows. *Tools and Trends: International Conference on Digital Preservation* Koninklijke Bibliotheek, 1-2 November 2007



Characterisation: General issues

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: General issues

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])



Categories of characteristics

Non-technical characteristics ("associated metadata", contextual)

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?

...



Categories of characteristics

Technical characteristics

Information concerning

- structure of the object
- internal representation of the object
- instructions on how to represent the object

e.g. dimensions of an image, image data,
compression, font information



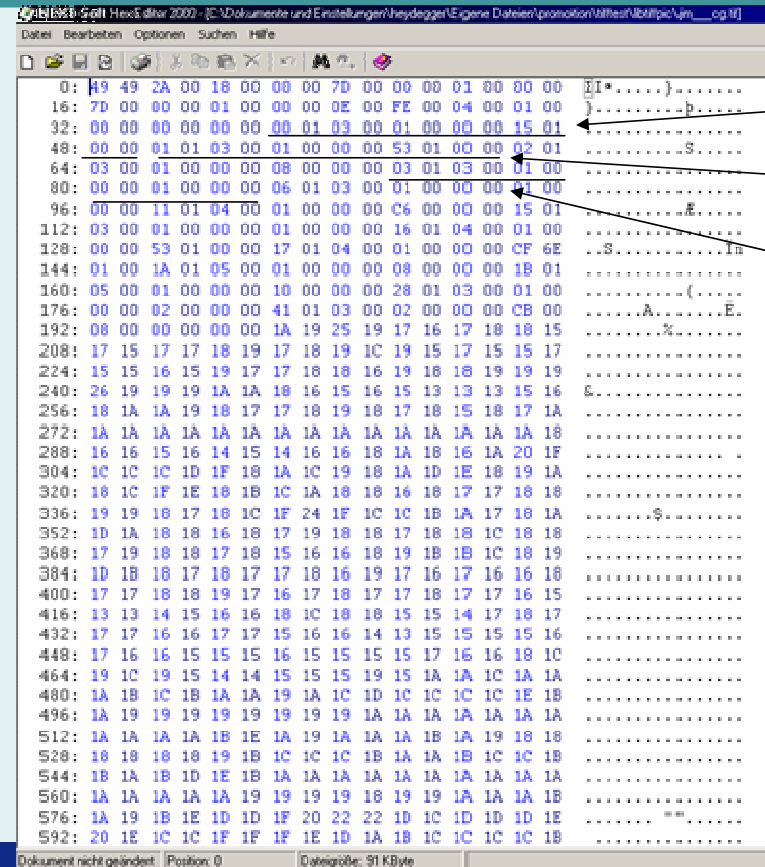
Characterisation: General issues

Where do we find characteristics? How are they described?

“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])





0: H9 49 2A 00 18 00 00 00 7D 00 00 00 01 80 80 00
16: 7D 00 00 00 01 00 00 00 0E 00 FE 00 04 00 01 00
32: 00 00 00 00 00 00 00 01 03 00 01 00 00 15 01
48: 00 00 01 01 03 00 01 00 00 00 53 01 00 00 02 01
64: 03 00 01 00 00 00 08 00 00 00 03 01 03 00 01 00
80: 00 00 01 00 00 00 06 01 03 00 01 00 00 00 01 00
96: 00 00 11 01 04 00 01 00 00 00 C6 00 00 00 15 01
112: 03 00 01 00 00 00 01 00 00 00 16 01 04 00 01 00
128: 00 00 53 01 00 00 17 01 04 00 01 00 00 00 CF 6E
144: 01 00 1A 01 05 00 01 00 00 00 00 00 00 00 1B 01
160: 05 00 01 00 00 00 10 00 00 00 28 01 03 00 01 00
176: 00 00 02 00 00 00 41 01 03 00 02 00 00 00 CB 00
192: 08 00 00 00 00 00 1A 19 25 19 17 16 17 18 18 15
208: 17 15 17 17 18 19 17 18 19 1C 19 15 17 15 15 17
224: 15 15 16 15 19 17 17 18 18 16 19 18 18 19 19 19
240: 26 19 19 19 1A 1A 18 16 15 16 15 13 13 13 15 16
256: 18 1A 1A 19 18 17 17 18 19 18 17 18 15 18 17 1A
272: 1A 1A 1A 1A 1A 1A 1A 1A 1A 1A 1A 1A 1A 1A 1A 18
288: 16 16 15 16 14 15 14 16 16 18 1A 18 16 1A 20 1F
304: 1C 1C 1C 1D 1F 18 1A 1C 19 18 1A 1D 1E 18 19 1A
320: 18 1C 1F 1E 18 1B 1C 1A 18 18 16 18 17 17 18 18
336: 19 19 18 17 18 1C 1F 24 1F 1C 1C 1B 1A 17 18 1A
352: 1D 1A 18 18 16 18 17 19 18 18 17 18 18 1C 18 18
368: 17 19 18 18 17 18 15 16 16 18 19 1B 1B 1C 18 19
384: 1D 1B 18 17 18 17 17 18 16 19 17 16 17 16 16 18
400: 17 17 18 18 19 17 16 17 18 17 17 18 17 17 16 15
416: 13 13 14 15 16 16 18 1C 18 18 15 15 14 17 18 17
432: 17 17 16 16 17 17 15 16 16 14 13 15 15 15 15 16
448: 17 16 16 15 15 15 16 15 15 15 15 17 16 16 18 1C
464: 19 1C 19 15 14 14 15 15 15 19 15 1A 1A 1C 1A 1A
480: 1A 1B 1C 1B 1A 1A 19 1A 1C 1D 1C 1C 1C 1E 1B
496: 1A 19 19 19 19 19 19 19 1A 1A 1A 1A 1A 1A 1A
512: 1A 1A 1A 1A 1B 1E 1A 19 1A 1A 1B 1A 19 18 18
528: 18 18 18 18 19 1B 1C 1C 1C 1B 1A 1A 1B 1C 1C 1B
544: 1B 1A 1B 1D 1E 1B 1A 1A 1A 1A 1A 1A 1A 1A 1A
560: 1A 1A 1A 1A 1A 19 19 19 18 19 19 1A 1A 1A 1B
576: 1A 19 1B 1E 1D 1D 1F 20 22 22 1D 1C 1D 1D 1E
592: 20 1E 1C 1C 1F 1F 1F 1E 1D 1A 1B 1C 1C 1C 1B

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.

Some facts about file formats

How many file formats are there?

- 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”

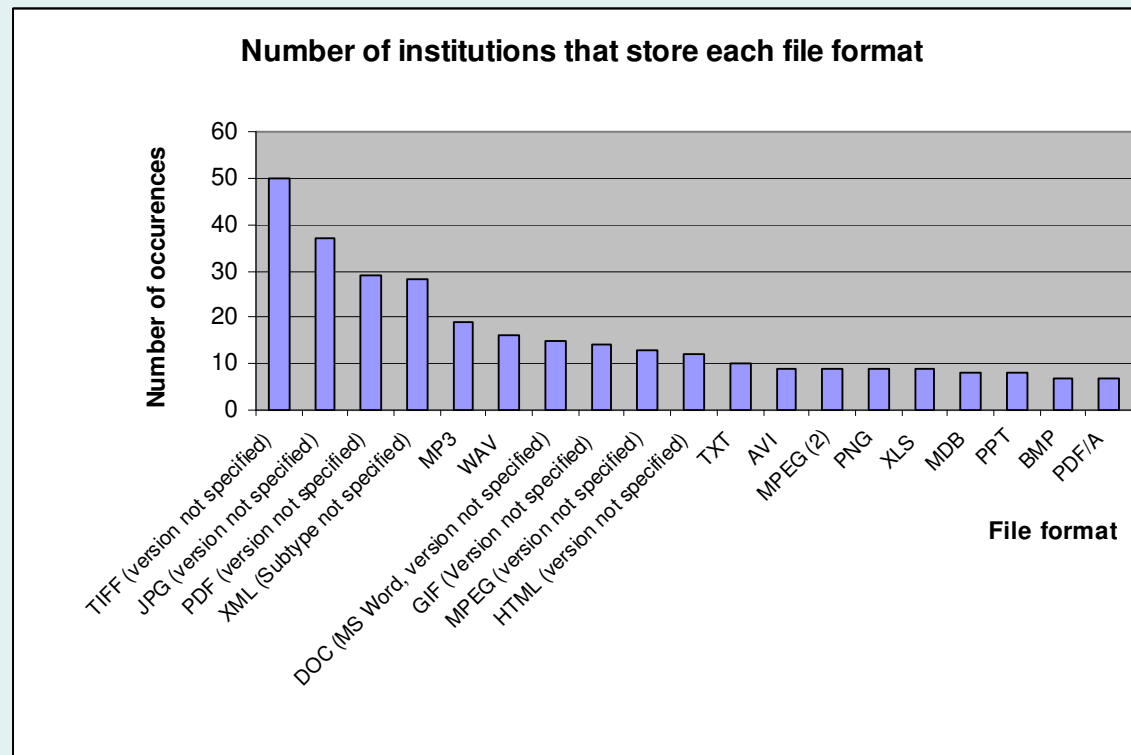
- 76 institutions from 13 countries
- 124 different file formats



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



How many file formats are used more often?



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



Suitability of formats for preservation (1)

High confidence	Medium confidence	Low confidence
<ul style="list-style-type: none">❖ TIFF (uncompressed)❖ PNG (*.png)	<ul style="list-style-type: none">❖ BMP (*.bmp)❖ JPEG/JFIF (*.jpg)❖ JPEG2000 (prefer lossless or uncompressed) (*.jp2)❖ TIFF (compressed)❖ GIF (*.gif)	<ul style="list-style-type: none">❖ MrSID (*.sid)❖ TIFF (in Planar format)❖ FlashPix (*.fpx)❖ PhotoShop (*.psd)❖ All other raster image formats not listed here



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



Suitability of formats for preservation (2)

High confidence	Medium confidence	Low confidence
<ul style="list-style-type: none">❖ Plain text (encoding: ISO8859-1 - 9, UTF-8, UTF-16 with BOM)❖ XML (includes XSD/XSL/XHTML, etc.; with included or accessible schema and character encoding explicitly specified)❖ PDF/A-1 (ISO 19005-1)	<ul style="list-style-type: none">❖ Cascading Style Sheets (*.css)❖ DTD (*.dtd)❖ PDF (*.pdf) (embedded fonts)❖ Rich Text Format 1.x (*.rtf)❖ HTML 4.x (include a DOCTYPE declaration)❖ SGML (*.sgml)❖ Open Office (*.sxw/*.odt)❖ Office Open XML (*.docx)	<ul style="list-style-type: none">❖ PDF (*.pdf) (encrypted)❖ Microsoft Word (*.doc)❖ WordPerfect (*.wpd)❖ DVI (*.dvi)❖ All other text formats not listed here



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



Suitability of formats for preservation (3)

High confidence	Medium confidence	Low confidence
<ul style="list-style-type: none">❖ AIFF (PCM) (*.aif, *.aiff)❖ WAV (PCM) (*.wav)	<ul style="list-style-type: none">❖ SUN Audio (uncompressed) (*.au)❖ Standard MIDI (*.mid, *.midi)❖ Ogg Vorbis (*.ogg)❖ Free Lossless Audio Codec (*.flac)❖ Advance Audio Coding (*.mp4, *.m4a, *.aac)❖ MP3 (MPEG-1/2, Layer 3)(*.mp3)	<ul style="list-style-type: none">❖ AIFC (compressed) (*.aifc)❖ NeXT SND (*.snd)❖ RealNetworks 'Real Audio, (*.ra, *.rm, *.ram)❖ Windows Media Audio (*.wma)❖ WAV (compressed) (*.wav)❖ All other audio formats not listed here



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



Criteria for suitability

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

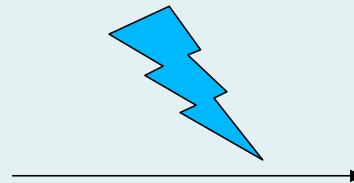
(J. Rog, C. van Wijk: Evaluating File Formats for Long-term Preservation, iPres 2007)



Example: Robustness of Formats

Robustness

::= resilience of file formats against bit-stream corruption



See: 'Just one Bit in a Million: On the Effects of
Data Corruption in Files'. <http://planets-project.eu/publications/?l=j>



In a nutshell...

- Characterisation is an essential part within an overall preservation framework.
- File Format is the central concept for representation of digital content, i.e. the characteristics of objects.
- There is a large number of formats which characterise objects in its own different way, but only a couple of them are actually suitable for preservation.

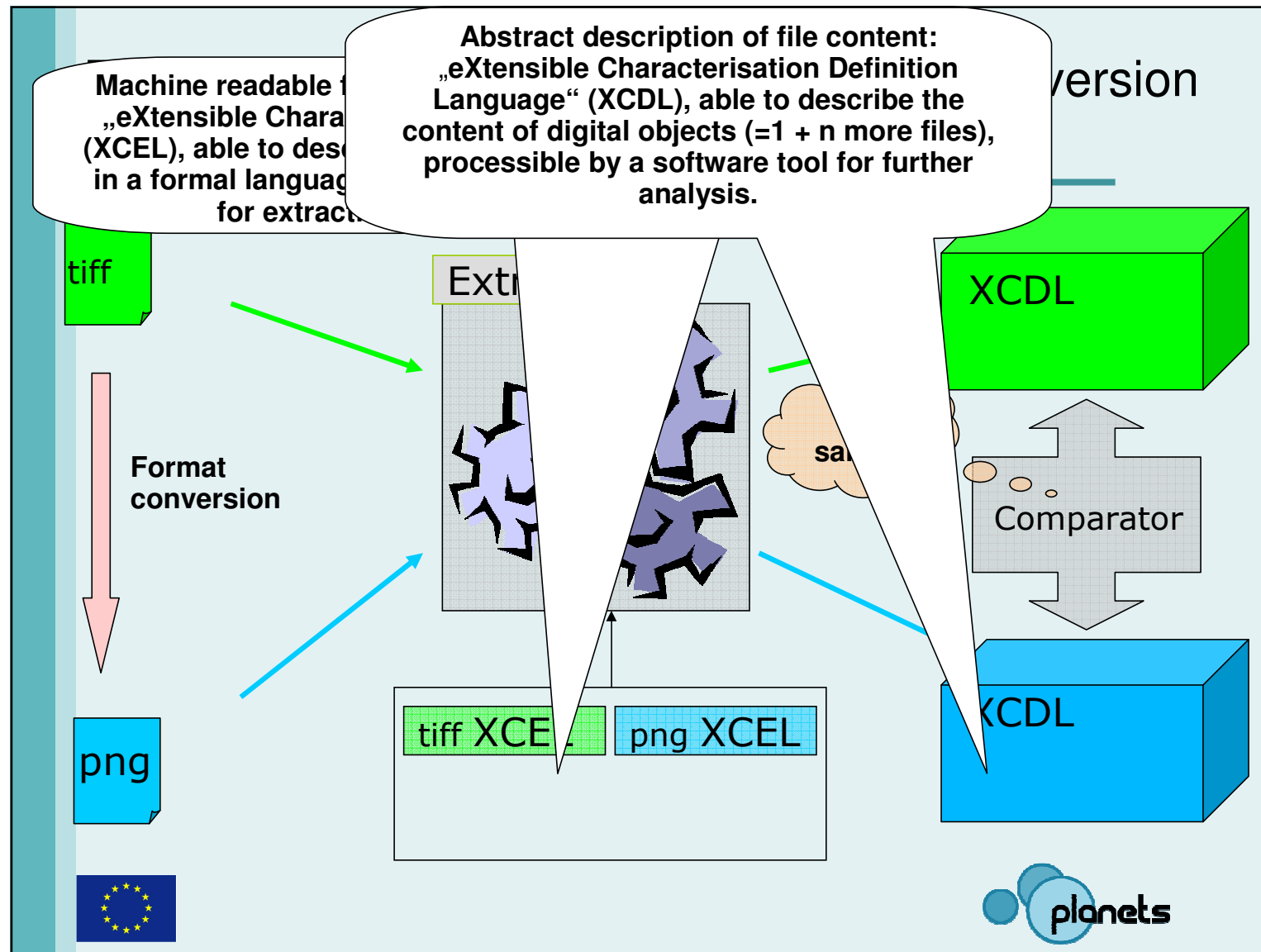


3

XCL overview

- Support preservation planning task
- 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





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: Basic Structuring Elements

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

processing

nonValidValues

valueInterpretation

item

param

valueLabel

symbol

value



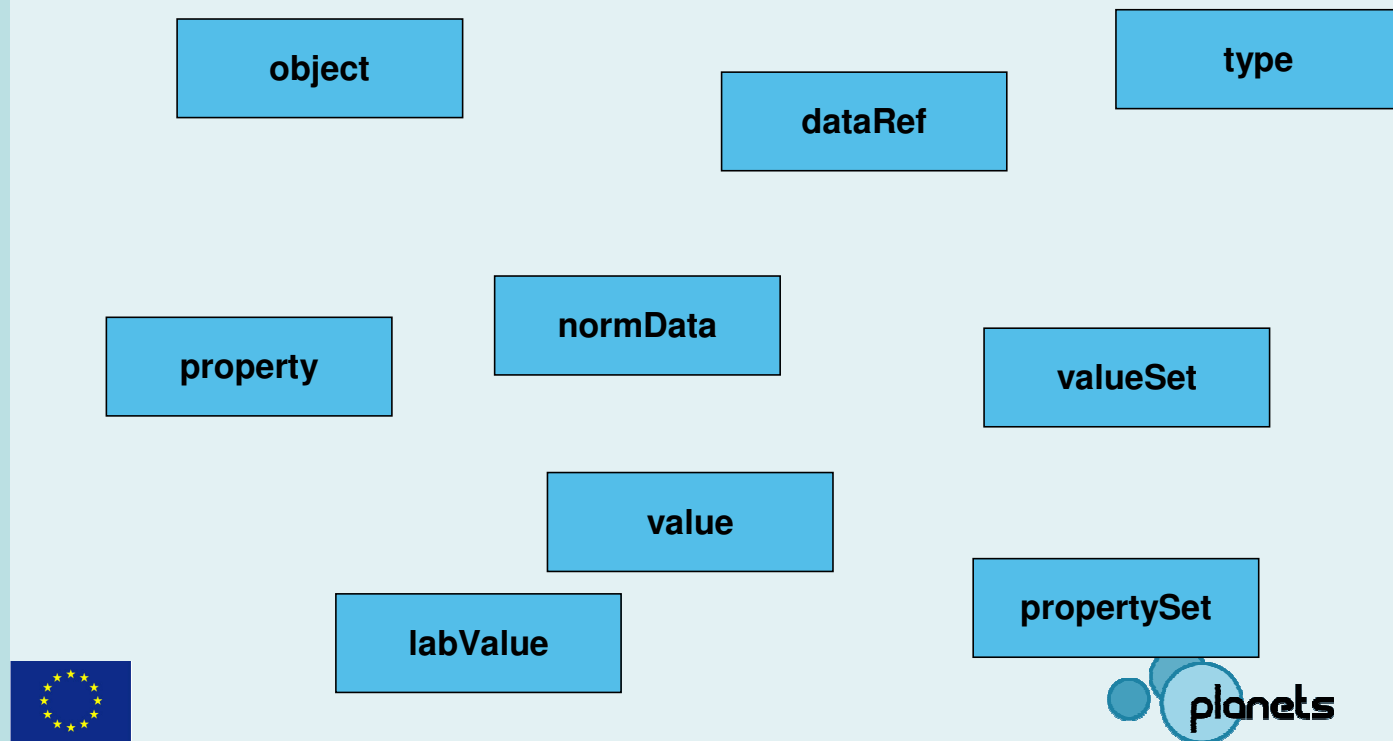
eXtensible Characterisation Definition Language (XC^{DL})

- Describes the content of a file /set of files in an abstract way.
- Designed for description 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:



Benefits of the XCL approach

- XCL is a generic solution, uses an abstract model, provides a unique vocabulary (‘XCL ontology’).
- 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



5

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

```
<!-- 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"/>
    [...]
```

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"/>
    [...]
```

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"/>
    [...]
```

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)
<pre><labValue> <val>32</val> <type>int</type> </labValue></pre>	<pre><labValue> <val>32</val> <type>int</type> </labValue></pre>
<p>copra output:</p> <pre>- <property id="2" name="imageHeight" unit="pixel" state="complete"> <metrics> <metric id="1" name="equal"> <result state="ok">true</result> </metric> </metrics> </property> -</pre>	



Thank you!

More about XCL on:

http://planetarium.hki.uni-koeln.de/planets_cms/index.php

